

# Mapping Streaming Architectures on Reconfigurable Platforms

Nikolaos Bellas, Sek M. Chai, Malcolm Dwyer, Dan Linzmeier  
 Embedded Systems Research, Motorola Labs  
 (*bellas@labs.motorola.com*)

**Abstract**— Hardware accelerators, used as application-specific extensions to the computational capabilities of a system, are efficient mechanisms to enhance the performance and reduce the power dissipation in a System On Chip (SoC). These accelerators execute on the computationally critical part of the application, and offload work from the main processors. In this paper, we present a design automation tool that generates accelerators based on a given application kernel. The accelerators are processing streaming data, and support a programming model which can naturally express a large number of embedded applications, and which results in efficient and fast hardware implementations. We demonstrate the applicability of the tool for architectural space exploration for a number of media applications, with results on area, throughput, and clock speeds.

## I. INTRODUCTION

The levels of integration of modern FPGAs have advanced to a point where the performance and flexibility are sufficient to map all functions of a complex SoCs into a single die. FPGA manufacturers have embedded fixed functionality cores such as general purpose processors, multipliers, multi-ported SRAM memories, and DSP slices in order to speed-up commonly used applications. At the same time, tool vendors have offered a plethora of pre-defined peripherals, fixed IP functions, and even synthesizable processor cores for the designer to customize the chip. The availability of a tool flow that abstracts out the particular hardware structures and presents a software-only front end interface to the application developer is a necessary step to precipitate the acceptance of FPGAs as SoC platforms.

Typically, scalar processors are reasonably efficient in handling normal conditional code with a low degree of instruction and data level parallelism. However, they are inefficient for high throughput, parallelizable code due to limited support of all kinds of parallelism (instruction, data, and task). They are further limited by the low memory bandwidth due to the narrow pipes to the main core.

In this paper, we describe an automation process which maps streaming data flow graphs (sDFG) to accelerators of the main scalar core. The streaming programming model assumes that the kernels process streams of data with a relatively limited lifetime, and deterministic memory access pattern. The streaming model decouples the description of memory access sequences from the computation within a kernel, thus making the customization of each of these two components (computation and memory access) easier and

more re-usable.

The design space exploration involves an iterative design cycle in which Pareto-optimal implementations of a given sDFG are produced under user and system constraints. For each iteration, a search space iterator instantiates a set of parameters that meet the given constraints, then a scheduler produces a schedule of operations optimized for throughput, and finally an RTL generation back-end tool produces the hardware description of the accelerator. Each generated accelerator is synthesized and implemented on an FPGA to be evaluated in terms of stream throughput, area, and clock speed, and later classified as Pareto-optimal or eliminated from consideration.

The contributions of the papers are the following:

- first, we propose the usage of the streaming paradigm (sDFGs) for application acceleration in a SoC-based reconfigurable fabric.
- second, we propose a template-based, automated method to perform architectural exploration on the accelerated application by evaluating the design space separately for the stream unit and the stream computational unit.
- and third, we explain how these concepts are placed in the context of a bus-based SoC design and how the accelerators are connected to the rest of the system.

The rest of the paper is organized as follows: Section II gives background information on the streaming programming paradigm and explains how it exploits technology trends that favor computation over communication. Section III details our template-based methodology, and section IV presents a set of embedded applications and the results of the method on a Xilinx Virtex-4 FPGA. Section V gives a summary of previous work on the relative areas, and Section VI presents the conclusion and future work.

## II. STREAM PROGRAMMING MODEL

### A. Architecture

The hardware accelerators generated by our method follow the streaming architectural paradigm [4]. They act as filters on input streaming data to generate the output streaming data specified by the streaming data flow graphs. Stream kernels exhibit a large degree of data and task level parallelism, with regular or even statically defined communication patterns [1].

The regularity of data access and the short lifetime of the

stream data allow for efficient optimization of the communication and the computational portions of the algorithm. Even more importantly, they make possible the decoupling of the stream access from the computation and their separate optimization.

Under this model, memory load/store operations no longer need to be scheduled amongst compute operations and optimal scheduling of operations now does not depend upon memory latencies. With this independence, the underlying memory system may be changed or may exhibit variable latencies, as with caches, with no effect on the computation schedule.

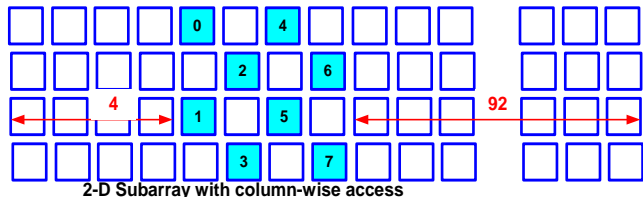
The decoupled memory access allows data pre-fetching to occur during computation. The programmer describes the shape and location of data in memory using stream descriptors. This decoupling allows the stream units to take advantage of available bandwidth to prefetch data before it is needed. The architecture becomes dependent on average bandwidth of the memory subsystem with less sensitivity to the peak latency to access a data element.

Deep pipelining allows multiple functional units to be chained, reducing access to large register files to store temporary data. This process is achieved with the programmer describing the data flow graph of the operations to be performed. Each operation is mapped to a set of functional units connected with a network. The number of functional units is dependent on the number of available logic gates, the number of potential parallel operations per cycle, and the user performance requirements.

### B. Stream Descriptors

The architecture includes several independent stream units to prefetch data from memory and turn streams into FIFO queues of stream elements. Additional stream units are created to write stream elements into memory. Each unit handles all issues regarding loading/storing of data including: address calculation, byte alignment, data ordering, and memory bus interface.

Data is transferred through the stream units which are programmed using stream descriptors. A stream descriptor is represented by the tuple (Type, Start\_Address, Stride, Span0, Skip0, Span1, Skip1, Size), where:



(Type, SA, Stride, Span[0], Skip[0], Span[1], Skip[1], Size) =  
(Byte, 4, 200, 2, -299, 2, -300, 8)

Figure 1. Stream descriptors for a memory access pattern

- Type indicates the element size in bytes (Type is 0 for bytes, 1 for 16-bit half-words, etc.).
- Start\_Address represents the memory address of the first stream element.
- Stride is the spacing, in number of elements, between two consecutive stream element.
- Span0 is the number of elements that are gathered before applying the skip0 offset.
- Skip0 is the offset applied between groups of span0 elements, after the stride has been applied
- Span1 is the number of elements that are gathered before applying the skip1 offset.
- Skip1 is the offset applied between groups of span1 elements, after the stride and the Skip0 have been applied.

The Stride, Span, Skip, and Type fields define the shape of the data stream in memory, while Start\_Address define the location of the first data element. Figure 1 shows an example of a static memory access pattern described by a two dimensional stream descriptor.

### C. Stream Computation

The streaming paradigm allows the application to exploit the large number of functional units that are readily available in modern VLSI technologies without taxing the communication resources [10]. We are using a “Data-flow Graph” (DFG) language to express operations in a machine-independent manner. A DFG consists of nodes, representing basic arithmetic, and logical operations composing the vector operation, and directed edges, representing the dependency of one operation on the output of a previous operation [7].

## III. TEMPLATE-BASED HARDWARE GENERATION

The problem we are addressing in this paper is the automatic generation of synthesizable accelerators from the streaming representation of Section II. Our approach is to select designs from a well-engineered framework, instead of generating the given hardware from a generic representation of a high level language. We generate highly optimized designs at various points at the cost-performance space based on the given application, the user requirements, and the capabilities of the rest of the system.

Figure 2 shows the iterative design flow. The main points of the tool flow are the following:

- a common template based on a regular architecture that accesses and processes streaming data,
- an iteration engine that instantiates system parameters that meet system and user constraints to initiate the next iteration of space search,
- a scheduler that performs sDFG scheduling and hardware allocation based on the parameters set by the iterator,
- an RTL constructor engine that produces optimized Verilog code for the data path and the stream units,

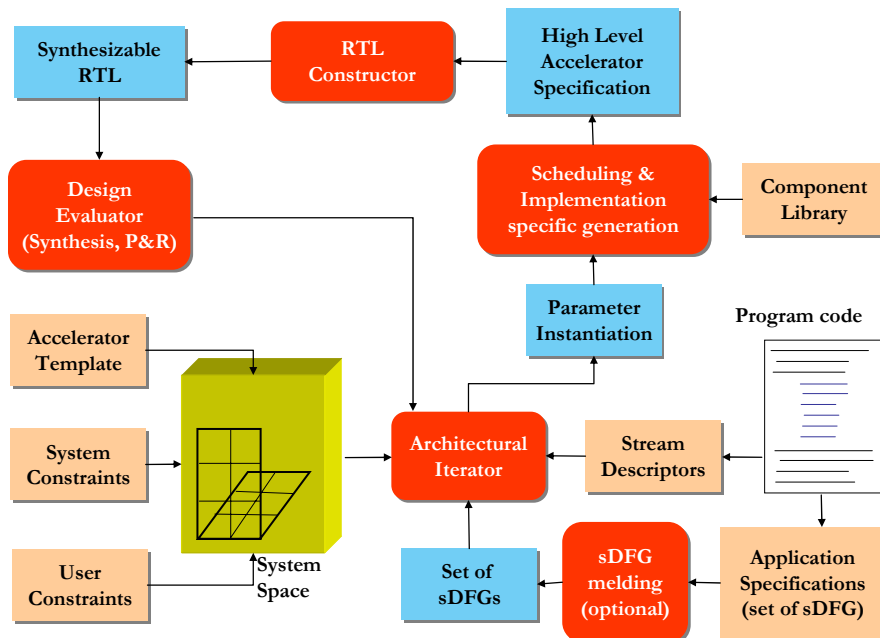


Figure 2 Template-based accelerator generation

- and an evaluation phase that synthesizes and maps the designs in FPGA and produces quality metrics such as area, and clock speed

Each of the data path and the stream unit have their own acceleration generation process. The rest of the section details each one of these engines and their interfaces.

#### A. Architectural template

The architectural template consists of two parts: the streaming data path and the stream unit (Figure 3). The stream unit expands into one or more input and output stream modules, and is generated to match the characteristics of the stream descriptors, and the characteristics of the bus-based system and the streaming data path. The data path is generated to execute a given sDFG to match user and system constraints in the specification space.

#### Stream Unit

The stream unit transfers streams from a system memory or peripheral, through a system bus and present them in-order to the accelerator. Likewise, it transfers processed output streams back to the memory.

The stream queue and the alignment unit store the incoming stream data and present them to the data path in-order. The number of storage elements, their size, and their interconnect depend on the stream descriptors and the requested bandwidth of the data.

The peak bandwidth for the accelerator depends on the schedule of the sDFG as we will discuss later. The size of the storage elements matches the size of the stream elements, for example it can be one byte for 8-bit pixel data.

Finally, the interconnect between the storage elements and the flow of streaming data between them depends on the span and skip of the stream description.

As we will examine later, the space iterator may also decide to allocate extra registers to the stream queue to match the system bus bandwidth capabilities. For example, in the case of an 8-byte system bus, the stream queue can have 8 or more storage elements to exploit the spatial locality of the memory accesses.

The bus line buffer is used to temporarily hold the data accessed from the system bus, and filter them to the stream queue when there is enough space. By detecting cases where the stride is greater than 1, the bus line buffer eliminates unnecessary elements before sending the stream to the stream queue.

The address generation unit (AGU) is hardwired to generate the memory access pattern of the stream descriptors. The number of registers that store internal variables, their width, the value and size of the stream description parameters are some of the configuration mechanisms of this unit.

The AGU aggressively generates addresses for data prefetching and sends them to the Address Line Buffer module. This module stores the addresses, merges addresses that fall within the same bus word (or burst size word in case bus burst is enabled), and competes for bus accesses with the other stream units. The generated number of buffers in the Address Line Buffer matches the average latency of the memory and bus systems and the capability of the bus to pipeline data accesses. For example, for a system bus that can pipeline up to two read accesses to a memory

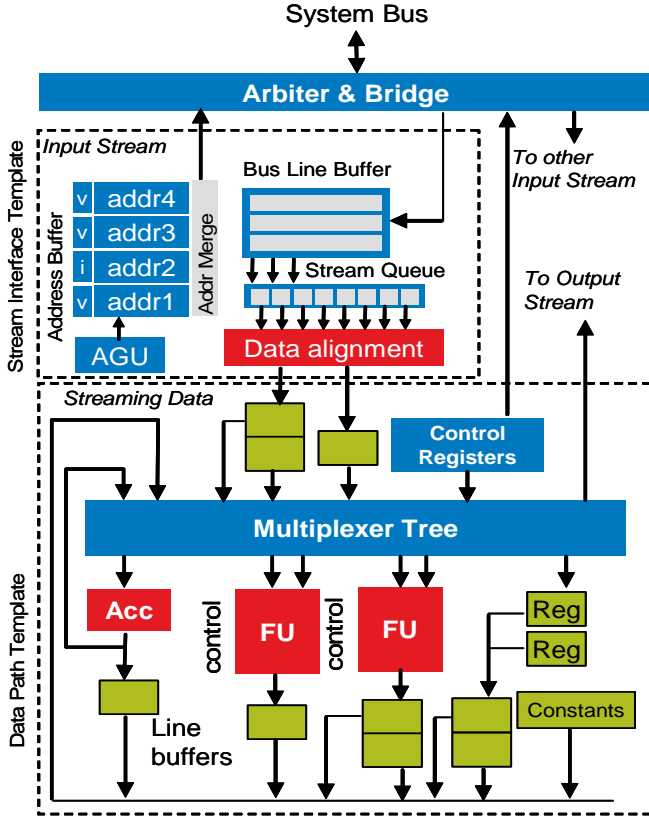


Figure 4. The accelerator template consists of the Data Path and the Stream Unit templates. Different optimizations are used in each case.

location the generator will create an Address Line Buffer with at least two address buffers.

Finally, the arbiter regulates the access of the stream units to the system bus. It uses a round-robin algorithm, and its complexity depends on the number of input and output streams of the sDFG.

#### Data path

The data path template of Figure 3 is an interconnect of reconfigurable functional units that produce and consume streaming data, and communicate via reconfigurable links. The links are chained at the output of a slice of a functional unit, and have a single input and potentially multiple outputs. They implement variable delay lines without the need of an explicitly addressable register file. The template also allow for the usage of a set of named registers that can be used by the sDFG to pass values from one sDFG iteration to the next and implement cross-iteration dependencies, and also to pass parameters to the program. Furthermore, the programming model allows for the use of accumulators for reduction operations.

The control logic of the data path is distributed and spatially close to the corresponding functional unit, multiplexer or line queue. This was an explicit design

decision to avoid creating critical paths due to long wires in modern VLSI technologies.

The type of the functional units (ALUs, multipliers, shifters, etc.), the specific operation performed within a type (e.g. only addition and subtraction for an ALU), the width of the functional unit, the size and number of storage elements of a FIFO, the interconnects between functional units (via FIFOs), the bandwidth from and towards the stream units are some of the reconfigurable parameters of the data path.

The data path requests data-sourcing from the input stream module and data-sinking from the output stream module. A simple, demand-driven protocol between the two modules is used to implement the communication. Stall signals from the stream units to the data path allow for a less than perfect memory system. A stall signal from any stream unit will cause the stall of the accelerator engine.

#### B. Architectural Iterator

The iterator selects a set of parameters in the space specified by the user and the system. For each one of this set of parameters, the tool flow builds an implementation by breaking the task into the implementation of the data path and the implementation of the stream unit.

#### Scheduling and High Level Implementation

The scheduler receives as input the sDFG along with the user and system constraints and schedules the operation of the sDFG to optimize throughput. The scheduler uses modulo scheduling to overlap multiple iterations in each cycle and exploits all the available parallelism under the resource constraints and data dependencies. The outline of the scheduling algorithm is given in Figure 3. The output of this stage is a hardware representation of the accelerator at a

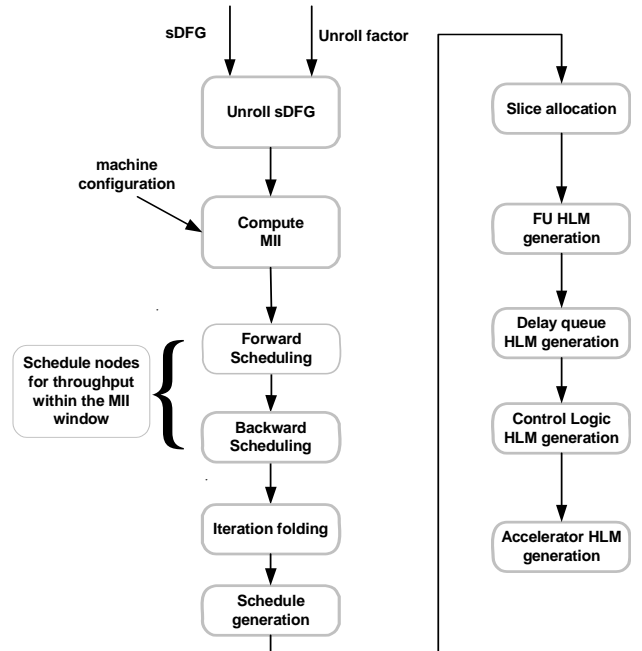


Figure 3. Scheduling and high level model (HLM) generation

higher specification level than an RTL specification (High level Model, HLM).

A strict lower bound of the initiation interval, called Minimum Initiation Interval (MII), is obtained by the number of available resources and the loop cross-iteration data dependencies [17].

During scheduling, the interconnects are not counted as resources, but rather they are “filled-in” during the generation of the HLM. By setting the schedule period equal to the MII, the scheduler maximizes the throughput of the accelerator, which is the main optimization target of the tool flow. Next, the schedule is generated within the MII window by first scheduling the nodes from top to bottom (forward scheduling) using a greedy approach. In this step, the nodes are scheduled immediately when all their parents have been scheduled and there exists an available resource to execute them.

Then, a backward scheduling heuristic is used to re-schedule some of the nodes by scheduling from the step MII-1 towards the step 0. This, in effect, “spreads out” the nodes within the steady-state period of the schedule and distributes the schedule more evenly within the MII steps. The net effect of this approach is to reduce the latency between successive nodes in the schedule, thereby reducing the storage requirements of the line delay queues [15].

The scheduler needs to only generate code for the steady state body of the schedule and not for the prologue and epilogue as is often the case in modulo scheduling [18]. Each data token that populates the FU inputs, outputs and line queues in every clock cycle is tagged with a valid bit. An operation produces valid output data only if both input data are valid. A source operation (like a stream load) produces data with valid bits when the data are available, and a sink operation (like a stream store) accepts data only when they are valid.

Next, the tool flow binds the operation nodes to the functional unit slices, and generates the delay links at the output of each slice to store the streaming outputs as they are produced by the FUs.

The stream unit design is generated based on user and system constraints. The size and number of buffer elements are chosen to meet the performance of the bus as well as the target performance of the generated data path. For example, the number of bus address queue elements, used to store pending addresses, is set to at least the bus pipeline factor so that bus transfers are sustained without stalling the data path. The number of line buffer elements, used to store data, should be at least the bus bandwidth to enable burst transfers. In addition, the number of stream data queue (used to store pending stream elements in a FIFO) is set to match the maximum bandwidth of the data path so that the stream unit can buffer the proper number of stream elements that can be consumed by the data path in a single cycle.

### C. RTL constructor

The RTL constructor reads the HLM representation and emits structural Verilog for the data path and the stream unit.

### D. Evaluator

At that point, the evaluation process is done by passing the resulting Verilog code through the Xilinx ISE tool-flow. We synthesize, and map the design targeting a Xilinx Virtex-4 architecture. We evaluate the design in terms of clock speed, and area overhead. As we will examine in the experimental evaluation, we are able to produce high-quality accelerators both in terms of area, and clock speed.

## IV. EXPERIMENTAL EVALUATION

### A. Methodology

This section describes the evaluation of the design methodology presented in previous sections. An application set, shown in Table 1, is selected from a wide range of media applications related to video compression, color processing, and image processing. Key compute intensive kernels from this application set is chosen for implementation. A design automation tool, using the design flow shown previously in Figure 2, is implemented in C++. The tool accepts each sDFG in the application set to generate candidate hardware accelerators according to the template shown in Figure 3. Different architectural configurations and loop unrolling factor are chosen such that Pareto-optimal designs for each benchmark can be chosen.

The generated hardware is synthesized and mapped onto a Xilinx Virtex4 LX60 FPGA, and the quality metrics of the produced bitstream (area, throughput, clock frequency) are recorded to assess the Pareto-optimality of the design

### B. Discussion

The results of Table 1 show the total number of FPGA, the average I/O bandwidth in bytes per cycle between the data path and the stream interfaces, and the clock frequency in

**Table 1 Representative performance and are results for a s et of multimedia benchmarks**

Streaming Kernels	FPGA Slices	Throughput (Bytes/cycle)	Frequency (MHz)
Binarization	189	1	218
Open Filter	731	0.15	188
Edge Detection	1677	0.07	174
Quantization	236	2	219
Column DCT	1148	1.23	171
Row DCT	1129	0.94	149
Color Processing (LPF)	2687	1.1	121
Color Processing (HPF)	2529	1.33	177

MHz after synthesis. These results enforce our initial premise that template-based approach can produce fast and area efficient designs. Using a high level representation such as the sDFG allowed for quick architectural exploration of different configurations. The streaming programming model facilitates the selection of sDFG selection and coding. Furthermore, it allows for design optimizations of both stream and data path, without recoding the benchmark.

In general, wider designs require more resources because the template design requires a larger number of queuing elements at the output of each functional unit to store live variables at each cycle. On the other hand, wider configurations are faster due to the lack of large multiplexers at the inputs of the functional units. In all cases, the maximum clock frequency is determined by the stream interface unit which is slower than the data path.

## V. RELATED WORK

In this section, we discuss previous work in the areas of streaming programming model and streaming architectures, architectural automation for ASIC and FPGA design flows, and special reconfigurable architectures.

### A. Streaming Programming Model and Architectures

Using streaming representations to expose concurrency and to express data communication explicitly has been recognized as an efficient way to both program off-the-shelf parallel processors, like graphics chips, and to architect new processors. A thorough analysis of the streaming programming model is given in [1]. Example streaming processors include Merrimac [10], Raw [20], Cell [14], and RSVP<sup>TM</sup> [7].

### B. Architectural automation for ASICs and FPGAs

There has been an intense interest in the research community in the last decade to automate the architectural process for ASIC of FPGA tool flows starting from a high level representation like C, Java, Matlab, DFGs and so on [9].

The PICO project incorporated a lot of concepts from earlier work on VLIW machines, and described a methodology to generate a VLIW engine along with an accelerator optimized for a particular application [18]. Similar projects include the Cyber tool [21], the OCAPI [19], and the DEFACTO compiler [23]. The Impulse-C [17] and Handel-C [24] tools are efforts to utilize C with extensions as a high level RTL language for FPGA design. At an even higher level of abstraction, the Matlab to gates compiler from AccelChip [3] targets mainly DSP kernels on FPGA platforms. Most of the above mentioned approaches use C as a more “user-friendly” hardware description language, and they add constructs to enhance concurrency, variable bitwidth, and so on in order to make C more amenable to hardware design. We believe that a template-based architectural automation that evaluates a large number of potential designs and focus on the most “profitable” parts

of the code is able to offer both design efficiency in terms of speed and cost, as well as programmability for developers that are not well-versed in hardware design.

A related problem is to automatically detect clusters of heavily executed assembly-level instructions that can be merged and extend the ISA of the processor. This research extends to both an ASIC environment [8] for the ARM processor.

### C. Reconfigurable processors

A number of academic projects and commercial products are tackling the hardware synthesis problem by designing efficient compile-time configurable or run-time reconfigurable architectures. This effort also stems from the fact that off-the-shelf, commercial FPGA architectures have little if any support for run-time reconfiguration. The GARP [5] and SCORE projects [6] propose the addition of reconfigurable planes that act as coprocessors of a scalar processor (MIPS in the case of Garp). Multiple planes offer a very fast context switch mechanism for run-time reconfiguration, and allows for virtual compute pages that can be mapped on the fabric both spatially and temporally.

Other projects include the RaPid architecture [11], the Chimaera architecture [22], the PipeRench architecture [13], and the RAW/Virtual Wires research [2].

## VI. CONCLUSION AND FUTURE WORK

A design methodology and prototype tool to automate the design and architectural exploration of hardware accelerators are described in this paper. These accelerators are programmed as streaming kernels to map to the streaming accelerators. In comparison to other approaches, we utilize a well-engineered template to enable fast convergence to an area and speed efficient design. We show how this methodology is used for an application set with various architectural configurations. New streaming accelerators are generated without recoding the application or re-design of the platform.

## REFERENCES

- [1] Amarasinghe S., Thies B. Architectures, Languages and Compilers for the Streaming Domain. Tutorial at the 12th Annual International Conference on Parallel Architectures and Compilation Techniques, New Orleans, LA
- [2] Babb J., et. al. Parallelizing Applications into Silicon. Proceedings of the 7th IEEE Symposium on Field Custom Computing Machines (FCCM), April 1999, Napa Valley, CA
- [3] Banerjee P. et. al.. A MATLAB compiler for distributed, heterogeneous, reconfigurable computing systems. Proceedings of the IEEE Symposium on Field Custom Computing Machines (FCCM), April 17-19, 2000, pp. 39-48, Napa Valley, CA
- [4] Nikolaos Bellas, Sek M. Chai, Malcolm Dwyer, Dan Linzmeier. FPGA implementation of a license plate recognition SoC using automatically generated streaming accelerators. 13<sup>th</sup> Reconfigurable Architectures Workshop (RAW), 25-26 April 2006, Rhodes, Greece
- [5] Callahan T., Hauser J., Wawrzynek J. The Garp Architecture and C Compiler. IEEE Computer Magazine, vol. 33, no. 4, April 2000, pp. 62-69
- [6] Caspi E., Huang R., Yeh J., Markovskiy Y., DeHon A., Wawrzynek J. Stream Computations organized for Reconfigurable Execution

- (SCORE): Introduction and Tutorial. BRASS research group technical report, University of California, Berkeley, August 2000
- [7] Chirisescu S., et. al. The Reconfigurable Streaming Vector Processor, RSVPTM. Proceedings of the 36th International Conference on Microarchitecture, December 2003, pp. 141-150, San Diego, CA
- [8] Clark N., Zhong H., and Mahlke S. Processor Acceleration Through Automated Instruction Set Customization. Proceedings of the 36th International Symposium on Microarchitecture, December 3-5, 2003, pp. 129-140, San Diego, CA
- [9] Compton K., Hauck S. Reconfigurable Computing: A Survey of Systems and Software. ACM Computing Surveys, vol. 34, No. 2, June 2002, pp. 171-210
- [10] Dally W. J., Hanrahan P., Erez M., Knight T. J., Labonté F., Ahn J.H., Jayasena N., Kapasi U. J., Das A., Gummaraju J., Buck, I. Merrimac: Supercomputing with Streams. Proceedings of the 2003 Supercomputing Conference, November 2003, pp-35-42, Phoenix, AZ
- [11] Ebeling C., Cronquist D., Franklin P., Secosky J., Berg S. Mapping Applications to the RaPiD configurable architecture. Proceedings of the 5th IEEE Symposium on Field Custom Computing Machines (FCCM), April 16-18, 1997, pp. 106-115, Napa Valley, CA
- [12] Gokhale M., Stone J., Arnold J., Kalinowski M. Stream-Oriented FPGA computing in the Streams-C High Level Language. Proceedings of the 8th IEEE Symposium on Field Custom Computing Machines (FCCM), April 17-19, 2000, pp. 39-48, Napa Valley, CA
- [13] Goldstein S. C. et. al. PipeRench: A Reconfigurable Architecture and Compiler. IEEE Computer Magazine, vol. 33, no. 4 April 2000, , pp. 70-77
- [14] Gschwind M., Hofstee P., Flachs B., Hopkins M., Watanabe Y., Yamazaki T. A novel SIMD architecture for the Cell heterogeneous chip-multiprocessors. Hot Chips XVII, August 15-16, 2005, Palo Alto, CA
- [15] Hwang C. T., Hsu Y. S., Lin Y. L. PLS: A Scheduler for Pipeline Synthesis. IEEE Transactions of Integrated Circuits and Systems, vol. 12, no. 9, September 1993, pp. 1279-1286
- [16] Kathail V., Aditya S., Schreiber R., Rau B.R., Cronquist D., Sivaraman M. PICO: Automatically Designing Custom Computers. IEEE Computer Magazine, vol. 35, no. 9, September 2002, pp. 39-47
- [17] Pellerin D., Thibault S. Practical FPGA Programming in C. Prentice Hall, 2005
- [18] Rau B. R. Iterative Modulo Scheduling. International Journal of Parallel Processing, 24:3-64, 1996
- [19] Schaumont P., Vernalde S., Rijnders L., Engels M., Bolsen I. A programming environment for the design of complex high speed ASICs. Proceedings of the 35th Design Automation Conference (DAC), June 1998, pp. 315-320, San Francisco, CA
- [20] Taylor M. B., et. al. The RAW Microprocessor: A Computational Fabric for Software Circuits and General Purpose Programs. IEEE Micro Magazine, 22(2), March 2002, pp.25-35
- [21] Wakabayashi K. and Okamoto T. C-based SoC design flow and EDA tools: An ASIC and system vendor perspective. IEEE Transactions on Computer-Aided Design, 19(12):1507-1522, December 2000
- [22] Ye A. Z., Moshovos A., Hauck S., Banerjee P. CHIMAERA: A high-performance architecture with a tightly-coupled reconfigurable unit. Proceedings of the 27th International Symposium on Computer Architecture (ISCA), June 2000, pp. 225-235, Vancouver, BC.
- [23] H. Ziegler H., Hall M. Evaluating Heuristics in Automatically Mapping Multi-Loop Applications to FPGAs Proceedings of the 13th International Symposium on FPGAs, February 2005, pp. 184-195, Monterey, CA
- [24] Celoxica Corporation, Handel-C language reference manual, [www.celoxica.com](http://www.celoxica.com)

RSVPTM is a trademark of Motorola Inc. Other product names are the property of their respective owner. A patent is pending that claims aspects of items and methods described in this paper.