

QuickDough: A Rapid FPGA Accelerator Design Method Using Soft Coarse-Grained Reconfigurable Array Overlay

Michael Shell, *Member, IEEE*, John Doe, *Fellow, OSA*, and Jane Doe, *Life Fellow, IEEE*

Abstract—FPGA accelerators used to offload compute intensive tasks of CPU offer both low power operation and great performance potential. However, they usually take much longer time to develop and are difficult to reuse, especially compared to the corresponding software design. Although the use of high-level synthesis (HLS) tools may partly alleviate this shortcoming, the lengthy low-level FPGA implementation process remains a major obstacle that limits the design productivity. To overcome this challenge, QuickDough, a rapid FPGA accelerator design method which utilizes soft coarse-grained reconfigurable arrays (SCGRAs) as an overlay on top of FPGA, is presented. Instead of compiling high-level applications directly as circuits implemented on the FPGA, the compilation process is reduced to an operation scheduling task targeting the SCGRA. Furthermore, the softness of the SCGRA allows domain-specific design of the processing elements, while allowing highly optimized SCGRA array be developed by a separate hardware design team. When compared to commercial high-level synthesis tools, QuickDough achieves competitive speedup in the end-to-end run time and reduces the compilation time by two orders of magnitude.

Index Terms—Overlay, FPGA Accelerator, Soft Coarse Grain Reconfigurable Array, Design Productivity

I. INTRODUCTION

The use of FPGAs as compute accelerators has been demonstrated by numerous researchers as an effective solution to meet the performance requirement across many application domains. However, despite years of research with numerous successful demonstrations, the use of FPGAs as computing devices remains a niche discipline that has yet to receive widespread adoption beyond highly skilled hardware engineers. Compared to a typical software development environment, developing an accelerator on FPGA is a lengthy and tedious process. While advancements in high level synthesis (HLS) tools have helped lower the barrier-to-entry for novel users, the irregularity of the synthesized circuit makes it difficult to compromise between execution time in cycles and implementation frequency. Furthermore, unless the HLS tool can directly synthesize the target FPGA configuration, vendor's back-end implementation tools must be employed for tasks such as floor planning and placing-and-routing. Compared to software compilation, the run-time of such back-end implementation tools is at least 2 to 3 orders of magnitude

slower, hindering the productivity of the designers especially during early development phases.

To address the above challenges, QuickDough, a rapid FPGA accelerator design method that utilizes SCGRA as an overlay, is proposed to produce FPGA bitstream directly from applications written in high-level languages such as C/C++. SCGRA has quite regular hardware structure and scales well on both the implementation frequency and execution time in cycles for parallel compute intensive kernels. Meanwhile, instead of being synthesized to circuits on the FPGAs, application compute kernels are translated to data flow graphs (DFGs) and further be *statically* scheduled to operate on the SCGRA. The lengthy hardware implementation tool flow is thus reduced to an operation scheduling problem.

Although the design and implementation of the SCGRA based accelerator must rely on the conventional hardware design flow, only one instance of the SCGRA design is required per application or application domain. Subsequent application development may then be accomplished rapidly by executing a simple scheduling algorithm. In addition, the performance of the SCGRA can be carefully optimized by a separate experienced hardware engineering team. Since the physical implementation of the SCGRA remains unchanged across applications or design iterations, the physical performance of the design can be guaranteed. According to the experiments on Zedboard, SCGRA based FPGA accelerator achieves competitive performance compared to Vivado HLS based accelerator under limited on chip buffer. While the SCGRA does consume more hardware overhead especially the BRAM blocks.

In Section III, the proposed FPGA acceleration design method QuickDough will be elaborated. The SCGRA implementation and compilation will be presented in Section V and Section IV respectively. Experimental results are shown in Section VI. Finally, we will discuss the limitations of current implementation in Section VII and conclude in Section VIII.

II. RELATED WORK

To improve the productivity of FPGA designers, researchers have approached the problem both by increasing the abstraction level and reducing the compilation time.

In the first case, decades of research in FPGA high-level synthesis have already demonstrated their indispensable role in promoting FPGA design productivity [12]. Numerous design languages and environments [10] have been developed to allow

M. Shell is with the Department of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA, 30332 USA e-mail: (see <http://www.michaelshell.org/contact.html>).

J. Doe and J. Doe are with Anonymous University.

Manuscript received April 19, 2005; revised January 11, 2007.

designers to focus on high-level functionality instead of low-level implementation details. While high-level abstraction may help a designers express the desired functionality, the low-level compilation time spent on synthesis, mapping, placing and routing for FPGAs remains a major hindrances to designs' productivity. Researchers have approached the problem from many angles, such as through the use of pre-compiled hard macros [21] in the tool flow, the use of a partial reconfiguration, and modular design flow [15].

On top of the above approaches, overlays, which can be parametric HDL Model, pre-synthesized or pre-implemented coarse-grained reconfigurable circuits over the fine-grained FPGA devices, promise both to raise the abstraction level and reduce the compilation time. Thus great research efforts have been attracted over the years and a number of overlays have been proposed [18, 19, 22, 23, 26, 28]. The granularity of these overlays ranges from multi-processors to highly configurable logic arrays.

Soft processors, which allow customization from various angles for target applications or application domains, have already been demonstrated to be efficient overlays for implementing an application on FPGA. [27],[2], and [3] employ general processors as the overlay and mainly have micro-architecture parameters such as pipeline depth configurable. [16] uses general processor with custom instruction extension as overlay, but hardware implementation is required whenever new custom instructions are added. [20] develops a fine-grain virtual FPGA overlay specially for custom instruction extension, which makes the custom instruction implementation portable and fast. [22] has customized data path on a many-core overlay, it could support both the coarse-grain multi-thread parallelism and data-flow style fine-grain threading parallelism. [26] adopts a multi-processor overlay with both micro-architecture and interconnection customizable. [23] and [28] develop reconfigurable vector processors as the FPGA overlay to cover larger domains of applications. [18] presented a GPU-Like overlay for portability and it could explore both the data-level parallelism and thread level parallelism. These processor level overlays typically approach the customization through instruction set extension or micro-architecture parameters tuning, and the application developers don't need much interaction with the low level hardware customization. Thus an application can be implemented rapidly, while the penalty is the hardware overhead and implementation frequency.

[8] and [17] build fine-grain components and mixed-grain components as a virtual FPGA overlay over the off-the-shelf FPGA devices. The virtual FPGAs allow the designers to reuse the virtual bitstream which is compatible across different FPGA vendors and parts. Particularly, the virtual FPGA with coarse granularity of components could also decrease the compilation time. [13] developed a family of intermediate fabrics which fits well for data parallel circuit implementation and the compilation time is almost comparable to software compilation. Apparently, the virtual FPGA overlays are beneficial to improving the design productivity and portability, though they do result in moderate hardware overhead and timing degradation.

Between the processor level overlays and virtual FPGA level

overlays, CGRA overlays on FPGA have unique advantages of compromising hardware implementation and performance especially for compute intensive applications as demonstrated by numerous ASIC CGRAs [25] [11]. CGRAs on FPGA and ASIC have many similarities in terms of the scheduling algorithm and array structure, however, they have quite different trade-off on configuration flexibility, overhead and performance. Basically, CGRAs on ASIC need to emphasize more on configuration capability to cover more applications, while FPGAs' inherent programmability greatly alleviate this concern. Accordingly, CGRAs on FPGA could accept more intensive customization while design productivity comes up as a new challenge.

[19] develops a VLIW architecture based parameterizable CGRA overlay called WPPA and provides an interconnection wrapper unit for each processing element to dynamically configure the topology of the CGRA. [14] proposes an heterogeneous CGRA overlay with multi-stage interconnection on FPGA, and the compilation can be done in milliseconds. While the CGRA size is quite limited and the implementation frequency is low due to the multi-stage interconnection. [24] develops a CGRA overlay named QUKU to improve reconfiguration speed for DSP algorithms. The experiments show that the CGRA overlay bridges the gap between soft processor and customized IP core. [9] builds an high speed mesh CGRA overlay using the elastic pipeline technique and achieves the maximum throughput. These CGRA overlay work typically take DFG as input and it is also possible to extract DFG during the software compilation. Thus the CGRA overlay helps to raise the abstraction level compared to conventional hardware design. Given the CGRA overlay architecture, the compilation doesn't involve any circuit optimization such as timing and pipelining at all and it makes the design accessible to an application developer without much hardware design experience. Especially, unlike the soft processor overlay and virtual FPGA overlay, [9] shows that the CGRA overlay doesn't have to compromise between the implementation frequency and design productivity as well as design portability.

As demonstrated in previous CGRA work, one of the major motivations of using CGRA overlays is its promising performance acceleration capability for compute intensive applications. However, a complete accelerator design on a hybrid general purpose processor (GPP) + FPGA using CGRA overlay is still missing. In fact, communication between FPGA and GPP is expensive, it has significant impact on the overall performance acceleration and thus influence the design choices of the CGRA overlay as well.

In this work, we opt to utilize a fully pipelined synchronous SCGRA as the overlay. Then we further implement it as an accelerator on Zedboard [5] which is a hybrid ARM + FPGA system. The accelerator now is configurable and is capable to handle all of our four full compute intensive applications with diverse data sets. With this SCGRA overlay based accelerator design method, an application can be implemented in a short time and the performance is competitive.

III. QUICKDOUGH FRAMEWORK

A. System Context

This work assumes a hybrid computing architecture with a host processor and an FPGA accelerator, where the processor handles tasks not-well suited to FPGA such as providing the OS environment and FPGA focuses on compute intensive kernels.

Figure 1 shows a typical FPGA acceleration architecture and all the experiments in this work stick to it. In this system, FPGA accelerator is attached to the system bus and it could access main memory through the bus. Inner the FPGA accelerator, there is a group of data buffers, an acceleration control block(Acc Ctrl), and a computation core. Data buffers are employed to store input data, output data and even temporary data of the computation. The Acc Ctrl block receives computation start signal probably from CPU via the bus when the input data is ready. Then it will trigger the computation core to start computation. When the computation is done, it will send the computation done signal back to interrupt CPU to collect data from the output buffer. The computation core is a customized circuit optimized for target compute kernel, which is supposed to be fast.

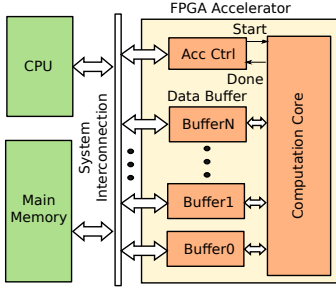


Fig. 1. A Typical FPGA Acceleration Architecture

B. QuickDough

In this section, QuickDough, which is developed to provide a customized acceleration solution to the target applications based on such a acceleration architecture, is presented. As shown in Figure 2, it starts from HW/SW partition where the compute intensive kernels of the applications are extracted. The extracted kernels are further transformed to DFGs, which are preferred in CGRA compilation. There are already work on both, which may help us to automate the processing [6] [4]. Currently, we just manually perform the HW/SW partition and DFG transformation in our preliminary research stage.

After the HW/SW partition, the design method can roughly be split into two parts. The part on top half is essentially the conventional software compilation except that we need to replace the compute kernel with accelerator drivers to manage the accelerator in the application, and the binary code is generated in the end. The part on the bottom half is basically the SCGRA compilation. Given the DFG of the compute kernel and accelerator configuration such as SCGRA size and on chip buffer capacity, the SCGRA compiler performs the operation scheduling, produces the configurations of the SCGRA overlay

and generates the bitstream eventually based on the pre-implemented SCGRA overlay in the SCGRA library. Finally, the binary code is loaded to the host system, the bitstream is configured to FPGA and the application is implemented on the hybrid compute system.

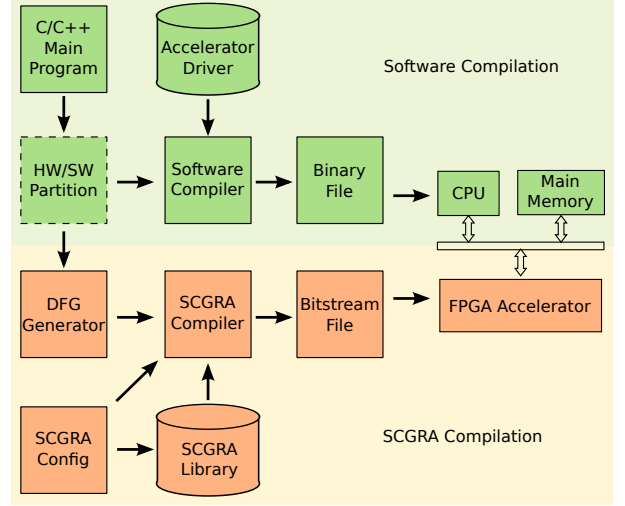


Fig. 2. QuickDough: An FPGA Accelerator Design Method Using SCGRA Overlay

IV. SCGRA OVERLAY INFRASTRUCTURE

One key idea of QuickDough is to rely on an intermediate SCGRA overlay to improve compilation time of the high-level user application. While the exact design of this SCGRA does not affect the compilation flow, its implementation does have a significant impact on the performance of the generated gateway.

A. SCGRA Based FPGA Accelerator

Figure 3 is the proposed FPGA accelerator built on top of the SCGRA overlay. The input/output data buffers and Acc Ctrl block are almost the same with those in the typical acceleration architecture 1, while the rest blocks are unique.

Particularly, the kernel part of the accelerator is a synchronous 2D torus SCGRA which consists of an array of PEs. Details of the PE will be illustrated in the next section. Another major difference is that two address buffers instead of the customized logic are used as the address generator to control the on-chip buffer accessing. With the address buffers, there is no need to develop any specific address generator when the target application changes, as we can simply replace its content together with the SCGRA configuration memory according to the scheduling result. Therefore, it reduces the chance of FPGA implementation and is beneficial to improving the design productivity.

Finally, it is noted that input and output data buffers could accommodate more data than that used by a single SCGRA execution, so the SCGRA may iterate multiple times before it consumes all the data in input buffer or fills the output buffer. From the perspective of host processor, it is more efficient to control multiple SCGRA execution at a time instead of

each SCGRA execution independently. In this case, we have a control unit called SCGRA Ctrl to make the multiple SCGRA execution transparent to the Acc Ctrl. Also sharing the data sets among multiple SCGRA execution could make best use of the limited input/output buffer. In addition, the regular SCGRA overlay tends to run at higher frequency than the data buffer controller handling the system bus protocol, and two sets of synchronous registers are added to keep the computation core and the rest of the design to running at individual clock domains.

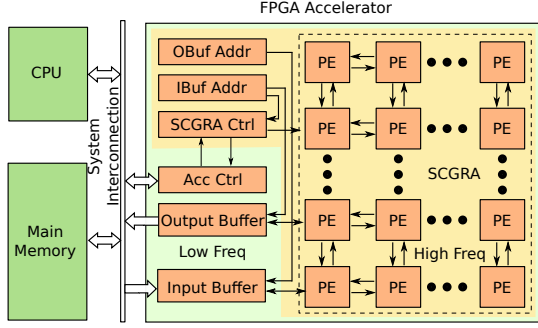


Fig. 3. SCGRA Accelerator

B. PE

In this section, an instance of PE is presented to demonstrate the feasibility of producing high performance gateway. As shown in Figure 4, the PE, centring an ALU block, a multiple-port data memory and an instruction memory, is highly optimized for FPGA implementation. In addition, load/store paths are implemented on the PEs that are responsible for data I/O beyond the FPGA. Addr Ctrl is used to start and reset the SCGRA execution by changing the instruction memory read address sequentially, so it has a single bit global start input signal from the SCGRA Ctrl block.

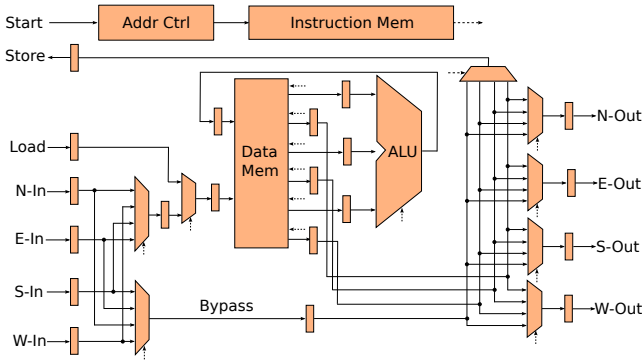


Fig. 4. PE structure

1) *Instruction Memory and Data Memory*: The instruction memory stores all the control words of the PE in each cycle. Since its content does not change during runtime, a ROM is used to implement this instruction memory and the content of the ROM can be loaded directly into a pre-implemented bitstream. The address of the instruction ROM is determined by the Addr Ctrl. Basically, the SCGRA execution will proceed

sequentially when the start signal is valid, and it will be reset when the start signal is invalid.

Data memory stores intermediate data that can either be forwarded to the PE downstream or be sent to the ALU for calculation. For fully parallelized operation, *at least* four read ports are needed – three for the ALU and one for data forwarding. Similarly, at least two write ports are needed to store input data from upstream memory and to store the result of the ALU in the same cycle. Although a pair of true dual port memories seem to meet this port requirement, conflicts may arise if the ALU needs to read the data while the data path needs to be written. As a result, a third dual port memory is replicated in the data memory.

Note that data memory here is usually implemented using a multiple port register file in many previous CGRA work. Although register file is even more flexible in terms of parallel reading and writing, the multi-port register file size is limited due to the inefficient hardware implementation. While we have a much larger DFG for scheduling and thus larger temporary storage is required, eventually a multi-port data memory is used instead.

2) *ALU*: At the heart of the proposed PE is an ALU designed to cover the computations in target application. Figure 5 shows an example design that could support all the operations of our benchmark which is listed in Table I. Complex operations like MULADD/MULSUB are implemented with DSP core directly. Operations with moderate complexity like ADDADD, RSFAND etc. are divided into two stages naturally and hardware block reuse is also considered at the same time. Finally, simple operations that can be done in a single cycle can be put in either stage depending on the pipeline status. Note that MUX in data path has significant influence on timing as well, so small MUX should be inserted properly and large MUX should be avoided. Currently, we just manually create the ALU design, but it is possible to automate this step.

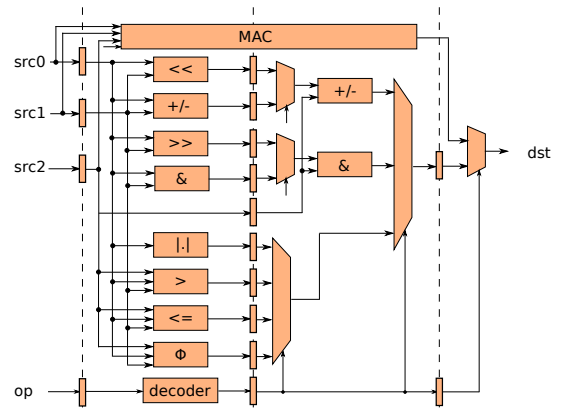


Fig. 5. ALU Example

C. Load/Store Interface

For the PEs that also serve as IO interface to the SCGRA, they have additional load path and store path as shown in 4. Load path and the SCGRA neighboring input share a single

TABLE I
OPERATION SET IMPLEMENTED IN ALU

Type	Opcode	Expression
MULADD	0001	$dst = src0 \times src1 + src2$
MULSUB	0010	$dst = src0 \times src1 - src2$
ADDADD	0011	$dst = src0 + src1 + src2$
ADDSUB	0100	$dst = src0 + src1 - src2$
SUBSUB	0101	$dst = src0 - src1 - src2$
PHI	0110	$dst = src0 ? src1 : src2$
RSFAND	0111	$dst = (src0 \gg src1) \& src2$
LSFADD	1000	$dst = (src0 \ll src1) + src2$
ABS	1001	$dst = abs(src0)$
GT	1010	$dst = (src0 > src1) ? 1 : 0$
LET	1011	$dst = (src0 \leq src1) ? 1 : 0$
ANDAND	1100	$dst = src0 \& src1 \& src2$

data memory write port, and an additional pipeline stage is added to keep the balance of the pipeline. Store path has an additional data MUX as well, but it doesn't have significant influence on the rest of the design.

V. SCGRA COMPILATION

Figure 6 illustrates the detailed SCGRA compilation of QuickDough. As shown in the diagram, the compilation essentially is to compile the compute kernel of an application to a specific SCGRA and generate the implementation bitstream.

The compilation starts from transforming the compute kernel probably written in high level language to a DFG as well as hyper block level IO mapping. When the DFG is determined, it can be scheduled to the specified SCGR using an operation scheduler. As the simulation performance of the DFG can be acquired from the scheduling, it is already enough to estimate the performance of the compute kernel. Moreover, taking the pre-built SCGRA implementation frequency and specified communication efficiency into consideration, it is possible to come up an even more accurate performance evaluation. Whether the performance design goal is met can be checked at this step. If no, we can go back to the DFG generation stage altering the DFG generation options such as loop unrolling factor. It may take multiple iterations to meet the design goal.

Once the design goal is achieved, the configuration words can be extracted from the scheduler and be further integrated into the pre-built SCGRA bitstream using the data2mem tool. Since the bitstream in the SCGRA library is bundled with specific FPGA device, HDL model will be used for porting to a new device and complete hardware implementation flow is required in that circumstance accordingly.

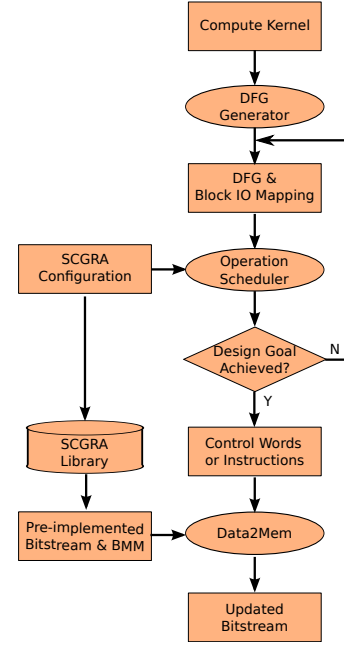


Fig. 6. SCGRA Compilation

A. DFG Generator

When the data set is small, it is usually possible to fully unroll the loops in the compute kernel and transform the unrolled kernel to DFG directly. However, when the data set gets larger, delicate unrolling strategy is needed as it has critical impact on both the performance and hardware overhead such as on chip buffer size and instruction memory depth. While unrolling strategy is not the focus of this paper, we will mainly explain the major constrains of the loop unrolling and the DFG generator. Also we will clear the required dumping list for the following compilation step.

Figure 7 shows how the DFG is extracted and executed on the SCGRA overlay based accelerator. The compute kernel can be performed by repeating the block execution, so the first step is to determine the block size that can be executed on the accelerator. Since the SCGRA overlay employs lock-step execution and the input/output data for each block execution must be fully buffered, the block size is mainly limited by the on chip buffer size. In this example, block size is set to be 10.

While the whole block can be completed by repeating the DFG execution, the next step is to decide the loop unrolling factor, such that the unrolled part can be transformed to DFG which can be executed on the SCGRA overlay at a time. In theory, the unrolling factor is probably limited by the instruction memory and data memory. Nevertheless, we have straightforward address buffers which store all the on chip buffer accessing addresses of the whole block execution. Although it is already set to be twice larger than the data buffer, it still overflows easily and becomes another major limitation of both the block size and the unrolling factor. In addition, the compute kernel depends on the repeating of the block execution and the block execution depends on the repeating of the DFG execution. As a result, the total loop iteration number must be fully divided by the block size. Similarly, the block size must also be fully divided by the

unrolling factor. This can be another unrolling and blocking limitation as well.

When the blocking and unrolling factor are determined, we can look at the dumping list of the DFG generator for the following compilation step. Apparently, DFG is required, as it is the input of the operation scheduler. And we store it in a text version as shown in Figure 7. The operation scheduler simply takes the DFG as input and it assumes the input/output of the DFG stored sequentially. However, it is more natural for the compute kernel to call the block computation directly. To bridge the gap, we have the block input/output stored sequentially, but have the address buffers to translate them back to the DFG preferred style. Therefore, the DFG generator must also provide the IO mapping information for each DFG execution as presented in Figure 7 as well.

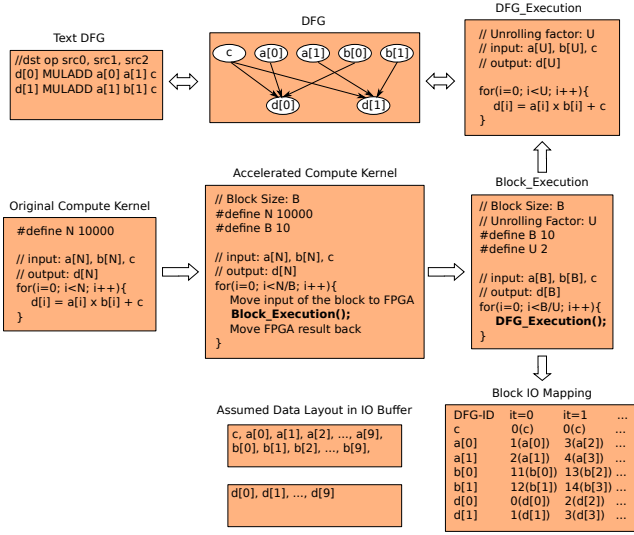


Fig. 7. DFG Generation

B. Operation Scheduler

In the SCGRA scheduling step, a list scheduling algorithm similar to [29] is adopted to schedule the DFG to the SCGRA infrastructure. It statically schedules the entire DFG on the target SCGRA. Such statically scheduled architecture is crucial in keeping the target SCGRA simple and efficient. To adapt to the proposed SCGRA structure, the scheduling metric is delicately adjusted to compromise the communication cost and load balance. On top of the DFG scheduling, it should also be responsible for generating the addresses for each SCGRA execution using the IO-mapper file.

C. Bitstream Integration

The final step of the compilation flow is to incorporate the instruction for each PE obtained from the scheduling stage with the pre-compiled SCGRA bitstream. By design, our SCGRA does not have mechanism to load instruction streams from external memory. Instead, we take advantage of the reconfigurability of SRAM based FPGAs and stored the cycle-by-cycle configuration words using on-chip ROM. The content of these instruction ROMs are embedded in the

configuration bitstream. In particular, the organization of the instruction ROM in the place-and-routed SCGRA design is obtained from its XDL file [7], which in turn allows us to create the corresponding BMM file. With this BMM file, the encoded instructions collected from the DFG scheduling may then be incorporated into the pre-implemented bitstream using the data2mem tool from Xilinx [1]. While original SCGRA design needs hours to implement, the bitstream updating scheme only costs a few seconds.

VI. EXPERIMENTS

In this work, we take four applications including matrix multiplication (MM), FIR filter, Kmean and Sobel edge detector as our benchmark. In order to be more representative, each application is further provided with three different data sets ranging from small, medium to large ones. The basic parameters and configurations of the benchmark is explained in the Table II.

The benchmark is implemented on Zedboard using both Vivado HLS and QuickDough to build the FPGA acceleration system. Then design productivity, implementation efficiency and performance of the two design methods are compared respectively.

A. Experiment Setup

All the runtimes were obtained from a computer with Intel(R) Core(TM) i5-3230M CPU and 8GB RAM. Vivado HLS 2013.3 is used to transform the compute kernel to hardware IP Catalog. Vivado 2013.3 is then used to integrate the IP core and build the FPGA accelerator on Zedboard. The SCGRA is initially built using ISE 14.7 and then exported as an IP core. Afterwards, it is integrated into the SCGRA based FPGA accelerator using PlanAhead 14.7. Finally, the accelerators are implemented on Zedboard to run the benchmark on bare-metal system mode.

Vivado HLS typically achieves the trade-off between hardware overhead and performance through altering the loop unrolling factor and loop blocking size. Larger block size is beneficial to data reuse and amortizing the communication cost, but it is usually constrained by the size of on-chip buffer. Similarly, larger loop unrolling factor typically promises better performance, however, hardware resources like DSP blocks will soon be used up and implementation frequency will degrade as well. In this work, two different sets of input/output buffer configurations are applied for the Vivado HLS based acceleration design. One set has fixed 2k-word buffer for both input and output, and the other set has maximum available buffer size and it could use as much as 64k-word for both input and output. Table III shows the near optimal blocking and loop unrolling setup for the benchmark.

SCGRA based FPGA accelerator has similar design choices to that of Vivado HLS in terms of loop unrolling and blocking. The difference is that the unrolled part of the loop is transformed to DFG which is further scheduled to the SCGRA infrastructure, and the loop unrolling factor is limited by the resources of the SCGRA overlay such as SCGRA size, instruction memory size and data memory size. While the

TABLE II
DETAILED CONFIGURATIONS OF THE BENCHMARK

Benchmark	MM	FIR	Sobel	Kmean
Parameters	Matrix Size	# of Input # of Taps+1	# of Vertical Pixels # of Horizontal Pixels	# of Nodes # of Centroids Dimension Size
Small	10	40/50	8/8	20/4/2
Medium	100	10000/50	128/128	5000/4/2
Large	1000	100000/50	1024/1024	50000/4/2

TABLE III
LOOP UNROLLING & BLOCKING SETUP FOR VIVADO HLS BASED ACCELERATOR DESIGN

Application		Max-Buffer		2k-Buffer		Complete Loop Structure
		Unrolling Factor	Block Structure	Unrolling Factor	Block Structure	
MM	Small	$2 \times 10 \times 10$	$10 \times 10 \times 10$	$2 \times 10 \times 10$	$10 \times 10 \times 10$	$10 \times 10 \times 10$
	Medium	$1 \times 1 \times 100$	$100 \times 100 \times 100$	1×100	10×100	$100 \times 100 \times 100$
	Large	1×500	50×1000	500	1000	$1000 \times 1000 \times 1000$
FIR	Small	2×50	40×50	2×50	40×50	40×50
	Medium	2×50	10000×50	2×50	1000×50	10000×50
	Large	2×50	50000×50	2×50	1000×50	100000×50
Sobel	Small	$1 \times 2 \times 3 \times 3$	$8 \times 8 \times 3 \times 3$	$1 \times 2 \times 3 \times 3$	$8 \times 8 \times 3 \times 3$	$8 \times 8 \times 3 \times 3$
	Medium	$1 \times 1 \times 3 \times 3$	$128 \times 128 \times 3 \times 3$	$1 \times 1 \times 3 \times 3$	$23 \times 128 \times 3 \times 3$	$128 \times 128 \times 3 \times 3$
	Large	$1 \times 1 \times 3 \times 3$	$75 \times 1024 \times 3 \times 3$	$1 \times 1 \times 3 \times 3$	$4 \times 1024 \times 3 \times 3$	$1024 \times 1024 \times 3 \times 3$
KMean	Small	$20 \times 4 \times 2$	$20 \times 4 \times 2$	$20 \times 4 \times 2$	$20 \times 4 \times 2$	$20 \times 4 \times 2$
	Medium	$5 \times 4 \times 2$	$5000 \times 4 \times 2$	$5 \times 4 \times 2$	$1000 \times 4 \times 2$	$1000 \times 4 \times 2$
	Large	$5 \times 4 \times 2$	$25000 \times 4 \times 2$	$5 \times 4 \times 2$	$1000 \times 4 \times 2$	$1000 \times 4 \times 2$

block size is constrained by both the size of input/output buffer and the addr buffer. Currently, we set the data buffer to be 2k-word and addr buffer to be 4k-word. However, addr buffer size is the major constrain of the blocking size as useless addresses are also stored. In this work, both a 2x2 SCGRA and 5x5 SCGRA overlay are used for all these benchmark implementation. Table IV and Table V present the detailed SCGRA overlay configuration and loop unrolling as well as blocking of the benchmark respectively.

TABLE IV
SCGRA CONFIGURATION

SCGRA Topology	Torus
SCGRA Size	$2 \times 2, 5 \times 5$
Data Width	32 bits
Instruction Length	72 bits
Instruction Memory Depth	1024
Data Memory Depth	256
Input/Output Buffer Depth	2048
Addr Buffer Depth	4096

B. Experiment Results

In this section, design productivity, hardware implementation efficiency and performance using both design methodologies are presented.

1) *Design Productivity*: Since design productivity involves many different aspects such as the abstraction level of the design entry, compilation time, design reuse, and design portability, it is difficult to evaluate all aspects especially some of the aspects are not easy to be quantized. In this section, we mainly look at the compilation time and design reuse which are critical angles of design productivity.

In order to accelerate an application using FPGA, Vivado HLS based design flow mainly consists of four steps including compute kernel synthesis, kernel IP generation, overall system

implementation and software compilation. At the compute kernel synthesis step, high level language program kernel is translated to a HDL model according to the user's synthesis pragma. Then the HDL model is further synthesized and packed to an IP core, which should meet the timing constrain and have the corresponding driver packed, at the kernel IP generation step. Afterwards, it is the system implementation step, where the IP core is integrated into the target system and the accelerator is implemented on FPGA. Finally, the application using the hardware accelerator is compiled to binary code. The SCGRA based FPGA accelerator design flow also needs four steps to complete the acceleration. The first step is translating the high level language program kernel to DFG, and the DFG is scheduled to a specified SCGRA overlay at the second step. At the third step, the bitstream is generated and corresponding driver is prepared. The last software compilation step is similar to that of the Vivado HLS based design flow.

Figure 8 and Figure 9 shows the compilation time needed to implement the benchmark using both Vivado HLS based design flow and the SCGRA based design flow respectively. Vivado HLS based design method needs lengthy IP core generation and hardware implementation for each application instance, which are the major time consuming processes. Compute kernel synthesis is usually as fast as the software compilation and can be done in a few seconds, but it may take up to 10 minutes when there is a pieplined large loop unrolling. Eventually, it takes the whole design flow 20 minutes to an hour to complete an implementation. With pre-implemented SCGRA overlay, the processing steps except the DFG scheduling of the SCGRA based design flow are quite fast and don't change much across the different applications. DFG scheduling is relatively slower especially when the DFG size and CGRA size are large, but it still can be done in a

TABLE V
LOOP UNROLLING AND BLOCKING SETUP FOR SCGRA BASED FPGA ACCELERATOR DESIGN

Application		SCGRA 2x2			SCGRA 5x5			Complete Loop Structure
		Unrolling Factor	DFG(OP/IO)	Block Structure	Unrolling Factor	DFG(OP/IO)	Block Structure	
MM	Small	$10 \times 10 \times 10$	1000/301	$10 \times 10 \times 10$	$10 \times 10 \times 10$	1000/301	$10 \times 10 \times 10$	$10 \times 10 \times 10$
	Medium	5×100	750/606	10×100	5×100	750/606	10×100	$100 \times 100 \times 100$
	Large	200	301/402	1000	200	301/402	10×100	$1000 \times 1000 \times 1000$
FIR	Small	40×50	860/131	40×50	40×50	860/131	40×50	40×50
	Medium	20×50	1000/141	100×50	50×50	2500/201	250×50	$10^4 \times 50$
	Large	20×50	1000/141	100×50	50×50	2500/201	250×50	$10^5 \times 50$
Sobel	Small	$4 \times 8 \times 3 \times 3$	1080/39	$8 \times 8 \times 3 \times 3$	$8 \times 8 \times 3 \times 3$	2160/55	$8 \times 8 \times 3 \times 3$	$8 \times 8 \times 3 \times 3$
	Medium	$4 \times 8 \times 3 \times 3$	1080/39	$8 \times 8 \times 3 \times 3$	$23 \times 8 \times 3 \times 3$	6210/115	$65 \times 8 \times 3 \times 3$	$128 \times 128 \times 3 \times 3$
	Large	$4 \times 4 \times 3 \times 3$	540/31	$16 \times 4 \times 3 \times 3$	$16 \times 4 \times 3 \times 3$	2160/55	$16 \times 4 \times 3 \times 3$	$1024 \times 1024 \times 3 \times 3$
KMean	Small	$20 \times 4 \times 2$	920/62	$20 \times 4 \times 2$	$20 \times 4 \times 2$	920/62	$20 \times 4 \times 2$	$20 \times 4 \times 2$
	Medium	$25 \times 4 \times 2$	1144/72	$125 \times 4 \times 2$	$125 \times 4 \times 2$	5768/272	$500 \times 4 \times 2$	$5000 \times 4 \times 2$
	Large	$25 \times 4 \times 2$	1144/72	$125 \times 4 \times 2$	$125 \times 4 \times 2$	5768/272	$500 \times 4 \times 2$	$50000 \times 4 \times 2$

few seconds. Finally, the SCGRA based design method could produce the bitstream in 5 to 15 seconds and it is typically two magnitude faster than the SCGRA overlay based design method.

expected, while the DSP consumption scales linearly with the SCGRA size. Note that the hardware overhead using Vivado HLS based design method in this table refers to the configuration with fixed 2k-word input/output buffer. The overhead with maximum input/output buffer configuration doesn't change the conclusion here and is skipped to save the space.

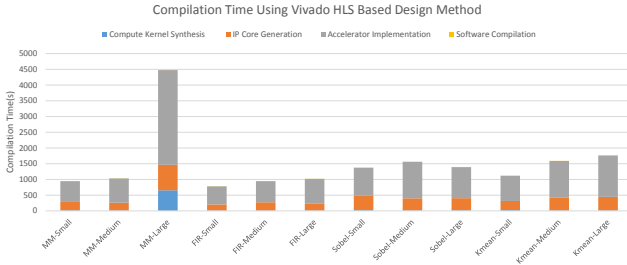


Fig. 8. Compilation Time Using Vivado HLS Based Design Method

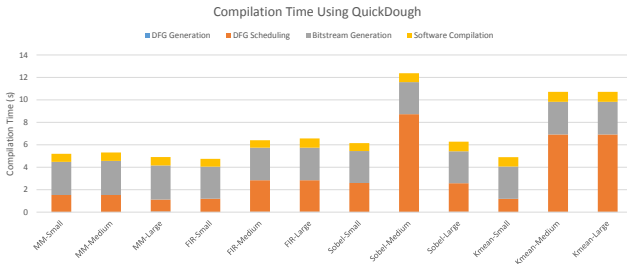


Fig. 9. Compilation Time Using SCGRA Overlay Based Design Method

2) *Hardware Implementation Efficiency*: In this section, hardware implementation efficiency including the hardware resource overhead, implementation frequency and implementation scalability using the two accelerator design methods are compared.

Table VI exhibits the hardware overhead using both accelerator design methods. It is clear that Vivado HLS based accelerator typically consumes much less FF, LUT and RAM36 due to the delicate customization for each application configuration. However, the number of DSP48 required increases dramatically with the expansion of the application kernel which is mostly induced by loop unrolling. As a result, DSP48 can be the bottleneck that constrains the performance of the accelerator through loop unrolling. SCGRA based accelerator usually costs more FF, LUT and particularly RAM36 as

TABLE VI
HARDWARE OVERHEAD USING BOTH VIVADO HLS AND THE SCGRA BASED ACCELERATOR DESIGN METHODS

		FF	LUT	RAM36	DSP48
MM	Small	5099	3730	4	102
	Medium	2778	3217	4	9
	Large	8282	10228	4	9
FIR	Small	5287	4494	4	99
	Medium	5283	4577	4	99
	Large	5283	4577	4	99
Sobel	Small	8503	6436	4	216
	Medium	6737	5526	4	144
	Large	6755	5546	4	144
KMean	Small	2852	3762	4	24
	Medium	5530	5427	4	24
	Large	5530	5427	4	24
SCGRA 2x2		9302	5745	32	4
SCGRA 5x5		34922	21436	137	25
FPGA Resource		106400	53200	140	220

Figure 10 shows detailed the implementation frequency of the benchmark using both FPGA accelerator design methods. Vivado HLS based accelerator works at most 100MHz on the Zedboard system due to the slow AXI controller. Even though the accelerator may work at individual clock domain, the synthesized IP core can also be slow especially when relatively larger loop unrolling is employed at high level language program synthesis step. While the SCGRA based is regular and it could works at higher clock domain. To make use of this advantage, the slow AXI controller block and the SCGRA are divided into two separate clock domains. Currently, SCGRA 2x2 works at 200MHz while SCGRA 5x5 works at 167MHz due to the high utilization of BRAM blocks. The controlling logic is slow and works at 100MHz.

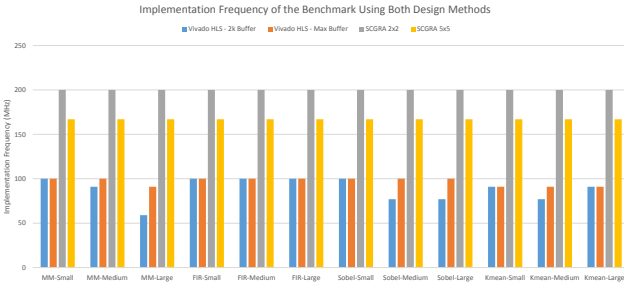


Fig. 10. Implementation Frequency of The Accelerators Using Vivado HLS and QuickDough

When the application data set gets larger, both the accelerator design methods may further scale the hardware infrastructure through loop unrolling and extending the SCGRA size etc. to accommodate the increased parallelism. Therefore, the scalability of the hardware implementation is also essential for the FPGA acceleration system. Figure xxx shows the implementation frequency of the accelerators kernel on xxx which could accommodate larger design than Zedboard. Since the implementation frequency of the Vivado HLS based accelerator depends on the specific application, we just use the simple matrix multiplication as an example.

3) *Performance*: In this section, the execution time of the benchmark is taken as the performance metric. Since the execution time of different applications and data sets varies a lot, the performance speedup relative to Vivado HLS based implementation with 2k-Buffer is used instead. Figure 11 shows the performance comparison of four different sets of accelerator implementations including two SCGRA overlay based accelerators and two Vivado HLS based accelerators. According to this figure, Vivado HLS based design method wins the MM-Medium, MM-Large, Sobel-Medium and Sobel-Large, while QuickDough outperforms in FIR with all three data sets, Kmean-Medium and Kmean-Large. The two design methods achieve similar performance on the rest of the benchmark.

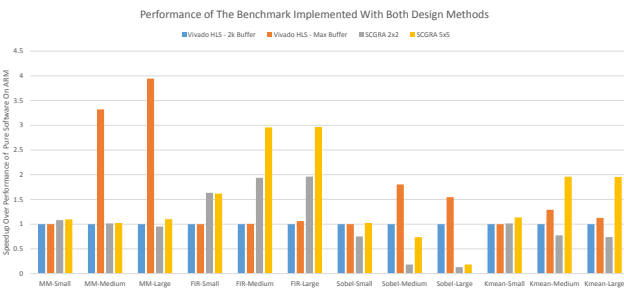


Fig. 11. Benchmark Performance Using Both Vivado HLS Based Design Method and The SCGRA Overlay Based Design Method

To insight in the comparison of the two design methods, Figure 12 presents the execution time decomposition of the benchmark. As the benchmark runs on an ARM + FPGA system where FPGA handles the compute kernel leaving the rest on ARM processor, the execution of the benchmark can be roughly divided into four parts: system initialization such as DMA initialization, communication time between FPGA

and ARM processor moving input/output data to/from the FPGA on-chip buffers, FPGA computation and the others such as input/output data reorganization for DMA transmission or corner case processing. The system initialization time is relatively small and is about the same for both design methods. The rest three parts determine the performance of different implementations. In addition, note that the execution time of different applications with diverse data sets varies in a large range. To fit all the data in a single figure, the execution time used in this figure is actually normalized to that of a basic software implementation on ARM.

According to this figure, Vivado HLS based designs especially the designs with max-buffer configuration perform better performance mainly due to the smaller overhead in communication and others which are essentially the input/output data reorganization time. If we further look at the blocking details in Table III, it is clear that Vivado HLS based designs with max-buffer configuration could accommodate both larger data sets and corresponding computation, which could improve data reuse and deduce the amount of communication as well as data reorganization cost eventually. However, when there is little data reuse between the neighboring blocks, the design with larger data block couldn't reduce the amount of communication. Moreover, as the DMA is used to transmit the data between FPGA and Main Memory and it brings additional cost such as DMA setup, communication cost per data can be reduced when the data block per DMA transmission is larger. The benefit will be negligible when the data set goes upto thousands and the additional DMA cost is amortized, while the exact number depends on the platform. In all the applications of the benchmark with medium and large data sets, Vivado HLS based designs especially the Max-Buffer configuration exhibit much smaller communication cost and data reorganization cost.

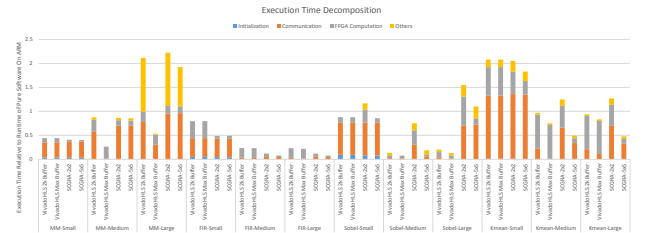


Fig. 12. Implementation Frequency of The Accelerators Using Vivado HLS and QuickDough

Computation time is one of the most significant part of the overall execution time. It basically depends on both the simulation performance in cycles and implementation frequency of the hardware infrastructure. As explained in previous sections, QuickDough using SCGRA overlay has regular structure and could provide scalable and higher frequency hardware infrastructure. Figure 13 shows the simulation performance which is the product of the block simulation performance and the number of blocks for each application. Apparently, QuickDough still fails to compete with the Vivado HLS based design in Sobel, but the performance gap of MM-Large is much smaller and particularly it shows better simulation

performance on almost all the rest of the benchmark. The simulation performance comparison is quite relevant to the depth of the loop unrolling as shown in Table III and Table V. It indicates that the simulation performance depends heavily on the depth of the loop unrolling no matter HLS based customized circuit or SCGRA based overlay is used as the hardware infrastructure. While QuickDough using SCGRA overlay as the hardware infrastructure has advantages in terms of implementation frequency, it enlarges its advantages in overall computation time as shown in Figure 14.

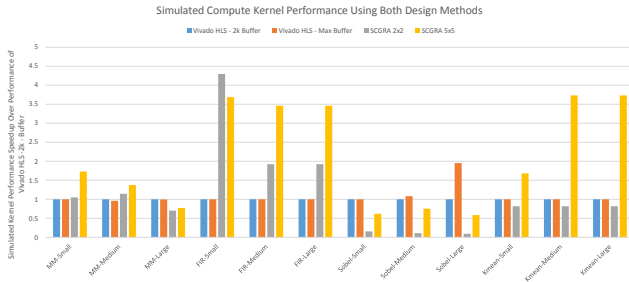


Fig. 13. Simulated Compute Kernel Performance Using Both Vivado HLS Based Design and QuickDough

In summary, Vivado HLS based design could afford larger buffer and accommodate larger block size, which further helps to reduce the communication time and the cost of input/output organization. Therefore, when there are more data reuse among neighboring blocks, the Vivado HLS based design is able to achieve better performance. QuickDough using SCGRA overlay could provide both higher simulation performance with larger loop unrolling in many cases and higher implementation frequency due to its regular structure, so it will win when the target application has smaller data set and more compute intensive kernels.

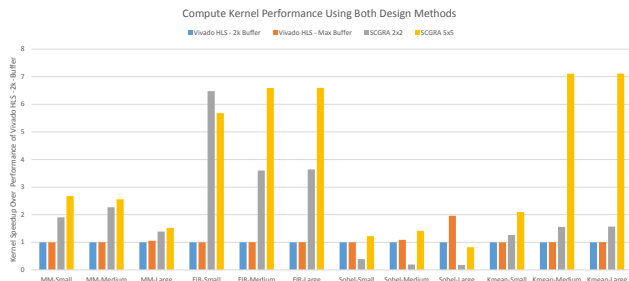


Fig. 14. Compute Kernel Performance Using Both Vivado HLS Based Design and QuickDough

VII. LIMITATIONS AND FUTURE WORK

While the current implementation of QuickDough has demonstrated promising initial results, there are a number of limitations that must be acknowledged and possibly addressed in future work.

First and foremost, the proposed methodology is designed to synthesize parallel computing kernels to execute on FPGAs only. As such, it is not a generic methodology to perform HLS on random logic. Furthermore, the proposed method is intended to serve as part of a larger HW/SW synthesis

framework that targets hybrid CPU-FPGA systems. Therefore, many high-level design decisions such as the identification of compute kernel to offload to FPGAs are not handled in this work. Also to guarantee the design productivity, a general front-end compilation that transforms high level language program kernel to DFG is still missing.

Currently, we just specify two SCGRA configurations for all the benchmark, while it is difficult for a high-level software designer to figure out an appropriate hardware configuration. An SCGRA optimizer will be developed to perform the SCGRA customization automatically in future.

Finally, the capacity of the address buffer used in QuickDough limits the block size that can be adapted to the FPGA in many cases. However, there are a large number of invalid address entries in it and this will be fixed in future.

VIII. CONCLUSIONS

In this paper, we have proposed QuickDough using a SCGRA overlay for compiling compute intensive applications to Zedboard. With the SCGRA overlay, the lengthy low-level implementation tool flow is reduced to a relatively rapid operation scheduling problem. The compilation time from an high level language application to the hybrid GPP+FPGA system is reduced by two magnitudes, which contributes directly into higher application designers' productivity.

Despite the use of an additional layer of SCGRA on the target FPGA, the overall application performance is not necessarily compromised. Implementation with higher clock frequency resulting from the highly regular structure of the SCGRA, in combination with an in-house scheduler that can effectively schedule operation to overlap with pipeline latencies provides competitive performance compared to a conventional HLS based design method.

REFERENCES

- [1] data2mem. http://www.xilinx.com/support/documentation/sw_manuals/xilinx11/data2mem.pdf. [Online; accessed 19-September-2012].
- [2] Microblaze soft processor core. <http://www.xilinx.com/tools/microblaze.htm>. [Online; accessed 25-June-2014].
- [3] Nios embedded processor. <http://www.altera.com/products/ip/processors/nios/nio-index.html>. [Online; accessed 25-June-2014].
- [4] Rocc2.0. <http://www.jacquardcomputing.com/roccc/>. [Online; accessed 19-January-2014].
- [5] Zedboard. <http://www.zedboard.org/>. [Online; accessed 25-June-2014].
- [6] Massimo Baleani, Frank Gennari, Yunjian Jiang, Yatish Patel, Robert K. Brayton, and Alberto Sangiovanni-Vincentelli. Hw/sw partitioning and code generation of embedded control applications on a reconfigurable architecture platform. In *Proceedings of the Tenth International Symposium on Hardware/Software Codesign, CODES '02*, pages 151–156, New York, NY, USA, 2002. ACM.
- [7] C. Beckhoff, D. Koch, and J. Torresen. The Xilinx design language (XDL): Tutorial and use cases. In

- Reconfigurable Communication-centric Systems-on-Chip (ReCoSoC)*, 2011 6th International Workshop on, pages 1–8. IEEE, 2011.
- [8] A. Brant and G.G.F. Lemieux. Zuma: An open fpga overlay architecture. In *Field-Programmable Custom Computing Machines (FCCM)*, 2012 IEEE 20th Annual International Symposium on, pages 93–96, 2012.
 - [9] D. Capalija and T.S. Abdelrahman. A high-performance overlay architecture for pipelined execution of data flow graphs. In *Field Programmable Logic and Applications (FPL)*, 2013 23rd International Conference on, pages 1–8, Sept 2013.
 - [10] J.M.P. Cardoso, P.C. Diniz, and M. Weinhardt. Compiling for reconfigurable computing: A survey. *ACM Computing Surveys (CSUR)*, 42(4):13, 2010.
 - [11] Katherine Compton and Scott Hauck. Reconfigurable computing: a survey of systems and software. *ACM Computing Surveys (csuR)*, 34(2):171–210, 2002.
 - [12] J. Cong, B. Liu, S. Neuendorffer, J. Noguera, K. Vissers, and Z. Zhang. High-level synthesis for FPGAs: From prototyping to deployment. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 30(4):473–491, 2011.
 - [13] J. Coole and G. Stitt. Intermediate fabrics: Virtual architectures for circuit portability and fast placement and routing. In *Hardware/Software Codesign and System Synthesis (CODES+ISSS)*, 2010 IEEE/ACM/IFIP International Conference on, pages 13–22, Oct 2010.
 - [14] R. Ferreira, J.G. Vendramini, L. Mucida, M.M. Pereira, and L. Carro. An fpga-based heterogeneous coarse-grained dynamically reconfigurable architecture. In *Proceedings of the 14th international conference on Compilers, architectures and synthesis for embedded systems*, pages 195–204. ACM, 2011.
 - [15] T. Frangieh, A. Chandrasekharan, S. Rajagopalan, Y. Iskander, S. Craven, and C. Patterson. PATIS: Using partial configuration to improve static FPGA design productivity. In *Parallel Distributed Processing, Workshops and Phd Forum (IPDPSW)*, 2010 IEEE International Symposium on, pages 1–8, april 2010.
 - [16] Mariusz Grad and Christian Plessl. Woolcano: An architecture and tool flow for dynamic instruction set extension on xilinx virtex-4 fx. In *ERSA*, pages 319–322, 2009.
 - [17] David Grant, Chris Wang, and Guy G.F. Lemieux. A cad framework for malibu: An fpga with time-multiplexed coarse-grained elements. In *Proceedings of the 19th ACM/SIGDA International Symposium on Field Programmable Gate Arrays, FPGA '11*, pages 123–132, New York, NY, USA, 2011. ACM.
 - [18] Jeffrey Kingyens and J. Gregory Steffan. The potential for a gpu-like overlay architecture for fpgas. *Int. J. Reconfig. Comp.*, 2011, 2011.
 - [19] Dmitrij Kissler, Frank Hannig, Alexey Kupriyanov, and Jürgen Teich. A dynamically reconfigurable weakly programmable processor array architecture template. In *ReCoSoC*, pages 31–37, 2006.
 - [20] D. Koch, C. Beckhoff, and G.G.F. Lemieux. An efficient fpga overlay for portable custom instruction set extensions. In *Field Programmable Logic and Applications (FPL)*, 2013 23rd International Conference on, pages 1–8, Sept 2013.
 - [21] C. Lavin, M. Padilla, J. Lamprecht, P. Lundrigan, B. Nelson, and B. Hutchings. HMFlow: Accelerating FPGA compilation with hard macros for rapid prototyping. In *Field-Programmable Custom Computing Machines (FCCM)*, 2011 IEEE 19th Annual International Symposium on, pages 117–124, may 2011.
 - [22] I. Lebedev, Shaoyi Cheng, A. Doupnik, J. Martin, C. Fletcher, D. Burke, Mingjie Lin, and J. Wawrzyniek. MARC: A many-core approach to reconfigurable computing. In *Reconfigurable Computing and FPGAs (ReConFig)*, 2010 International Conference on, pages 7–12, dec. 2010.
 - [23] A. Severance and G. Lemieux. Venice: A compact vector processor for fpga applications. In *Field-Programmable Technology (FPT)*, 2012 International Conference on, pages 261–268, Dec 2012.
 - [24] S. Shukla, N.W. Bergmann, and J. Becker. Quku: a two-level reconfigurable architecture. In *Emerging VLSI Technologies and Architectures*, 2006. IEEE Computer Society Annual Symposium on, pages 6 pp.–, March 2006.
 - [25] R. Tessier and W. Burleson. Reconfigurable computing for digital signal processing: A survey. *The Journal of VLSI Signal Processing*, 28(1):7–27, 2001.
 - [26] D. Unnikrishnan, Jia Zhao, and R. Tessier. Application specific customization and scalability of soft multiprocessors. In *Field Programmable Custom Computing Machines*, 2009. FCCM '09. 17th IEEE Symposium on, pages 123–130, April 2009.
 - [27] P. Yiannacouras, J.G. Steffan, and J. Rose. Exploration and customization of fpga-based soft processors. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 26(2):266–277, Feb 2007.
 - [28] Peter Yiannacouras, J. Gregory Steffan, and Jonathan Rose. Fine-grain performance scaling of soft vector processors. In *Proceedings of the 2009 International Conference on Compilers, Architecture, and Synthesis for Embedded Systems, CASES '09*, pages 97–106, New York, NY, USA, 2009. ACM.
 - [29] Hayden Kwok-Hay. Yu, Colin Lin. So. Energy-efficient dataflow computations on FPGAs using application-specific coarse-grain architecture synthesis. In *Highly Efficient Accelerators and Reconfigurable Technologies, The 4th International Workshop on*. IEEE, 2012.