

CHO: A Benchmark Suite for OpenCL-based FPGA Accelerators

University of Manchester Technical Report

GEOFFREY NDU, MIKEL LUJAN AND JAVIER NAVARIDAS

Manchester 2 May 2014

The Advanced Processors Technologies Research Group
School of Computer Science
The University of Manchester, Manchester, UK

UNIMAN-COMP-APT-TR-02-05-2014

CHO: A Benchmark Suite for OpenCL-based FPGA Accelerators

Geoffrey Ndu, Mikel Lujan, Javier Navaridas
School of Computer Science

University of Manchester, Manchester, M13 9PL, UK

Email: geoffrey.ndu, mikel.lujan, javier.navaridas@manchester.ac.uk



Abstract—Programming FPGAs with OpenCL-based high-level synthesis frameworks is gaining attention with a number of commercial and research frameworks announced. However, there are no benchmarks for evaluating these frameworks. To this end, we present CHO benchmark suite an extension of CHStone, a commonly used C-based high-level synthesis benchmark suite, for OpenCL. We characterise CHO at various levels and use it to investigate compiling non-trivial software to FPGAs.

1 INTRODUCTION

Open Computing Language (OpenCL) [1] is an open standard for platform-independent, general purpose parallel programming across CPUs, GPUs and accelerators. OpenCL consists of an API for coordinating parallel computation and a cross platform programming language (subset of ISO C99 with extensions for parallelism). It allows software to be written once and executed on any device that supports OpenCL. Its execution model consists of a host device which submits computational intensive kernels to compute devices for execution.

Programming Field-Programmable Gate Arrays (FPGAs) with OpenCL-based High-Level Synthesis (HLS) frameworks is now becoming mainstream with active support by the major FPGA vendors [2]. HLS is the automatic conversion of an algorithmic description into either a low-level Register Transfer Level (RTL) description or a digital circuit [3]. RTL refers to the low-level design abstraction that models a digital circuit in terms of the flow of digital signals between registers and the logical operations performed on those signals. HLS allows a designer to work more productively at a higher level of abstraction and achieve faster time-to-market than using error prone and difficult to debug RTL. Further, frameworks that use software programming languages such as C, OpenCL open up the power of FPGAs to software engineers (who outnumber hardware engineers by an order of magnitude).

Benchmarking is an important technique for analysing the performance of systems by studying the execution of the benchmark applications that are chosen to be a representation of the applications of interest. A good

HLS benchmark suite should allow HLS framework developers to qualitatively evaluate new ideas as well serve as a standard for benchmarking the diverse HLS frameworks available. From our discussion with HLS users (especially non-FPGA experts), who have to choose from the myriad of HLS frameworks, the second objective is equally important as the first.

In this paper, we introduce CHO: a suite of benchmark applications for OpenCL-based HLS platforms that meets the objectives set out above. CHO is an rewrite of the C-based CHStone benchmark suite [4]. CHStone is the commonly used HLS benchmark suite and consists of 12 applications from diverse application domains. Although based largely on C, OpenCL differs in some aspects from the standard C language. For example, OpenCL has disjoint memory spaces and moving data from one memory space to another need to be done explicitly. Hence, moving from one language to the other is often not straightforward.

This paper makes the following contributions:

- We present CHO an OpenCL port of the CHStone HLS benchmark allowing for the benchmarking of OpenCL-based HLS.
- We characterise CHO at various levels.
- We use CHO and a state-of-the-art OpenCL HLS framework to evaluate compiling non-trivial programs to FPGA.

2 RELATED WORK

Benchmarking of HLS frameworks does not have a rich history when compared to benchmarking of software platforms and compilers. Early HLS framework developers tend to use their own choice of applications for evaluation. In the 90's, the HLS community attempted to standardize benchmarking by proposing the 1992 High-Level Synthesis Workshop Benchmarks [5] and the 1995 High-Level Synthesis Design Repository [6]. These benchmarks covered a number different applications and application domains but were mostly written in algorithmic VHDL. VHDL is a type of Hardware Description Language

(HDL) (specialized language for encoding the structure, design and operation of electronic circuits). However, HLS frameworks has since moved from HDLs to high-level software languages, mostly variants of C. These benchmark suites have a few C-based applications but they are mostly tiny (less than a 100 lines of code) Digital Signal Processing (DSP) kernel loops [4]. Consequently, these benchmarks are rarely used nowadays.

CHStone [4] is now the de-facto standard benchmark suite used in the HLS community but it lacks support for OpenCL. There are OpenCL benchmark suites used in evaluating heterogeneous computing platforms such as Valar [7], Rodina [8] and SHOC [9] but they are too large and complex for FPGA synthesis and fitting.

Commercial and research OpenCL frameworks are increasingly being developed [2]. Notable frameworks include OpenCL-to-Silicon framework [10] and Altera SDK for OpenCL (AOCL) [11] (the first HLS to pass OpenCL conformance tests [11]). To the best of our knowledge none of the OpenCL frameworks were benchmarked with any standard benchmark suite.

3 BACKGROUND

3.1 FPGAs

An FPGA is essentially a sea of Look-Up Tables (LUTs). A LUT is a small high-speed memory and is programmed by loading a function's truth table as shown in Figure 1. Combining a LUT and a D flip-flop (a circuit with two stable states and can be used to store state information) results in what is often referred to as a logic cell. Several logic cells together with special-purpose circuitry such as an adder/subtractor carry chain form a logic block. Logic blocks can be connected to other logic blocks through a reconfigurable routing network making it possible to implement complex functions. FPGAs also contain other components such Block RAMs (BRAMs), DSP slices, various communication interfaces (e.g. PCI Express) and even processor cores. Figure 2 shows a typical high-end FPGA.

Modern FPGAs can have up to 2 million logic cells and 68 Mbits of BRAMs. Typically, FPGAs run at a much lower clock frequency than CPUs and GPUs but can outperform them by implementing custom and often more power efficient execution pipeline. [12], [13]. Figure 3 shows how a slower FPGA can outperform a CPU. The figure compares the implementation of 32-bit integer bit reversal [14] on an ARM 9 processor that requires about 38 instructions and 48 cycles, to an FPGA where the same operation can be performed by simply reversing the connections between two buffers. It is assumed that CPU and the FPGA implementation operate at the same frequency.

Despite the advantages offered by FPGAs their use is still largely restricted to a narrow segment of hardware programmers as programming often involves writing complex RTL code. HLS raises the level of abstraction (and productivity) by allowing algorithmic descriptions

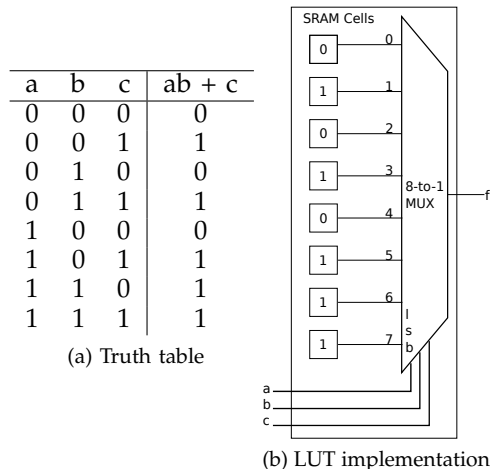


Fig. 1. Implementing logic function $f(a, b, c) = ab + c$

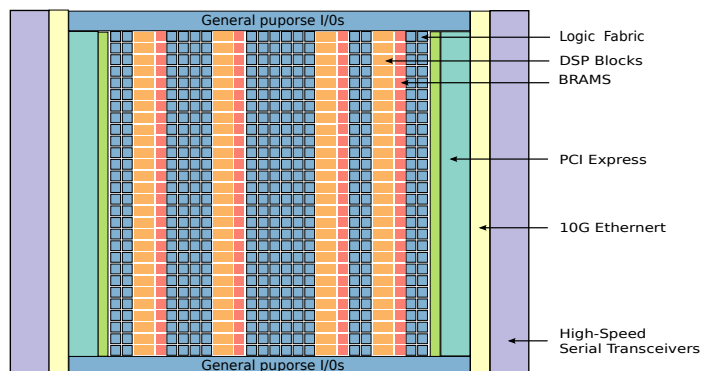


Fig. 2. FPGA architecture

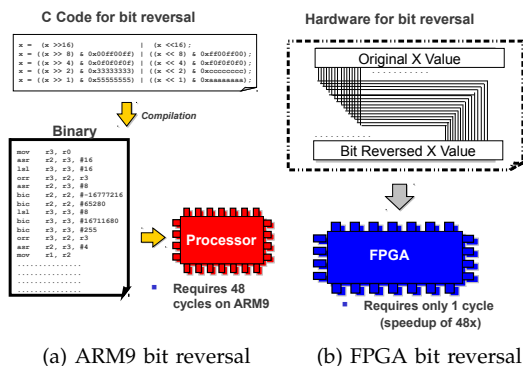


Fig. 3. How FPGAs outperform CPUs and GPUs

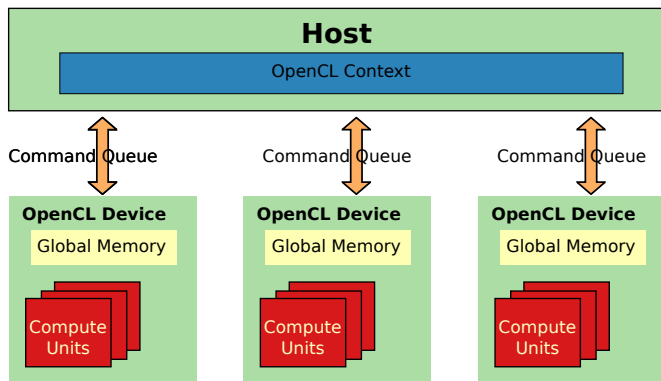


Fig. 4. OpenCL Platform Architecture

to be converted into RTL or even a digital circuit. Most HLS frameworks are now C-based and often support only a subset of the parent language.

3.2 OpenCL Architecture

An OpenCL computing platform consists of a CPU host connected to one or more OpenCL devices as shown in Figure 4. The part of an application targeting the devices is called the kernel. A kernel is hardware agnostic and should run on any device that supports the OpenCL standard. Kernels are written in OpenCL C programming language, a variant/subset of C-99.

The other part of an OpenCL application, the host, executes on the CPU and submits commands to perform the computations expressed in the kernels on a device. A kernel instance is called a work-item and multiple instances can be grouped into independent work-groups.

Two types of programming models are explicitly supported within OpenCL; the data parallel programming model and the task parallel programming model. In the task parallel programming model, a kernel is executed using a single work-item within a work-group while in the data parallel programming model multiple work-items are employed with input data partitioned across work-items.

The host defines a context for the execution of the kernels. The context includes the following resources: devices (devices to be used by the host), kernels (functions to run on devices), program objects (program sources and executables that implement the kernels) and memory objects (memory objects visible to the host and devices which contain values that can be operated on by the kernel instances). The context is created and manipulated by the host via OpenCL API.

Work-items have access to 4 disjoint memory regions: global memory (read/write memory accessible by all work-items), constant memory (read-only memory accessible to all work-items), local memory (read/write memory local to a work-group) and private memory (read/write memory private to a work-item)

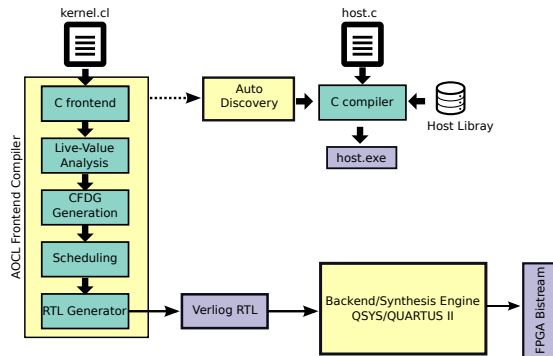


Fig. 5. Altera OpenCL SDK (AOCL)

3.3 OpenCL HLS Compilers

A typical OpenCL HLS framework can be divided into two parts, the front-end and the back-end. The front-end converts the kernels to RTL code while the back-end compiles the RTL down to FPGA bitstream (the configuration data that is loaded into the FPGA).

Figure 5 shows the structure of AOCL [15], the first commercial OpenCL-based HLS framework. The front-end leverages LLVM compiler infrastructure [16]. LLVM is increasingly used as the front-end of C-based HLS frameworks. The C-language front-end parses the kernel into LLVM Intermediate Representation (IR). LLVM IR uses simple RISC-like instructions augmented with high-level information such as types and explicit control flow graphs to represent programs (see Figure 6 and Figure 7). Notice that the LLVM instructions in Figure 7 are simple enough to translate directly into hardware operations (e.g., a load from memory, or an arithmetic computation). Hence, the rest of the front-end operates directly on the IR [15]. Each basic block is analysed and a control data flow graph is created for operations within a block. Each basic block is a hardware module that takes inputs from either the kernel arguments or another basic block, processes the data according to the instructions within and then produces output that is passed to other basic blocks. The basic-blocks are connected together to produce a complete circuit.

The scheduler attempts to minimize execution time as well as area. Finally, the RTL generator converts the IR instructions into Verlog RTL. The compiler automatically generates interfaces for accessing the off-chip memory (which serves as OpenCL's global memory) and the PCI Express (used by the host program to access the global memory). The traditional hardware compilation tools (synthesis tools) are used to synthesize the RTL code for FPGA.

The host program can be compiled with any C/C++ compiler and linked against AOCL's host library. This library implements the OpenCL function calls that launch the kernels on an FPGA. A module embedded inside the FPGA, termed the Auto-Discovery module, allows the host program to query the FPGA about the kernels it

holds. AOCL performs offline compilation only.

AOCL creates pipelined circuits to increase performance and maps work-items to pipeline stages. Let us assume that AOCL has created a three stage (load, add and store) pipelined hardware for the kernel in Figure 6, as shown in Figure 8. On the first clock cycle, work-item 0 is clocked into the load stage. The circuit begins by fetching the first elements of data from arrays a and b. On the second cycle, work-item 1 is clocked in while work-item 0, which has completed its read from memory and stored the results in the registers, moves to the second stage of the pipeline. Processing completes when the last work-item (work-item 7) exits from the last stage of the pipeline. For work-groups with a single work-items (i.e. kernels written in a task parallel fashion), AOCL would attempt to pipeline loops within kernels, using loop pipelining [17], to improve performance.

```
__kernel void vector_add(__global int *a, __global int *b, __global int *c)
{
    int i = get_global_id(0);
    c[i] = a[i] + b[i];
}
```

Fig. 6. Vector addition in OpenCL

```
entry:
%call = tail call @i32_to_i32_bitcast @get_global_id @i32_to_i32 @i32_0 #2
%idxprom = sext i32 %call to i64
%arrayidx = getelementptr inbounds i32* %a, i64 %idxprom
%0 = load i32* %arrayidx, align 4, !tbaa !2
%arrayidx1 = getelementptr inbounds i32* %b, i64 1
%1 = load i32* %arrayidx1, align 4, !tbaa !2
%add = add nsw i32 %1, %0
%arrayidx3 = getelementptr inbounds i32* %c, i64 %idxprom
store i32 %add, i32* %arrayidx3, align 4, !tbaa !2
ret void
```

Fig. 7. Vector addition in LLVM IR

4 CHO BENCHMARK SUITE

4.1 Overview

CHO extends the popular C-based CHStone HLS benchmark suite providing a means of benchmarking present and emerging OpenCL frameworks. The kernels in CHO are implemented using OpenCL task parallel model (i.e. one work-item per work-group) as it requires minimal change to the overall structure of the original code. Further, some of the programs are difficult to express using the data parallel programming model.

CHO, as its progenitor, consists of 12 diverse applications as shown in Table 1. CHO targets OpenCL 1.0, the minimum standard, allowing it to be compiled by any OpenCL compliant compiler as OpenCL is designed to be backward compatible. Like CHStone we don't prescribe a particular way of using CHO.

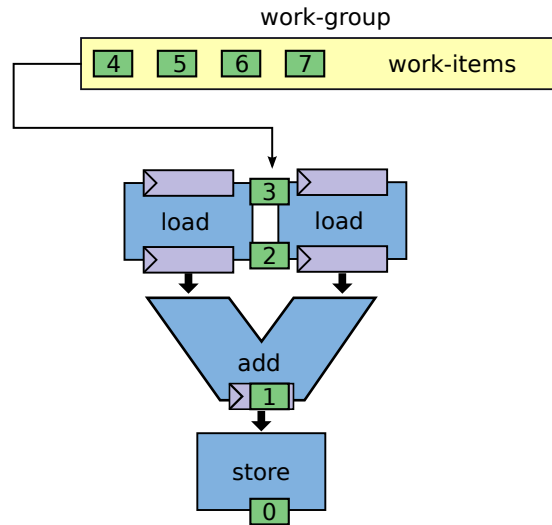


Fig. 8. AOCL pipeline implementation

Porting CHStone to OpenCL required eliminating the use of non-constant global variables as global variables in OpenCL must be constants. We made extensive use of structures (`struct`) to pass state from one function to another. To improve performance we allocated variables whenever possible in the private memory space as it is by definition the fastest memory in OpenCL. To allow compilers to perform memory optimization we added `restrict` keyword to pointer arguments whenever possible. The `restrict` informs compilers that pointers that share the same type do not alias.

Each application in CHO is a single kernel and has a test-bench that is part of the host program. Each kernel reads/writes its input/output data from/to the host side (via PCIe and the global memory). Implementing the test-bench on the host instead of encoding in the kernel enables the study of how an OpenCL HLS compiler would handle external I/O interfacing between the FPGA and the host (assuming that is done automatically by the framework). I/O interfacing and its bandwidth and latency are critical components in FPGA design. Often implementations need to be specialized based on achievable external input and output bandwidth. We reused the test-vectors from CHStone.

4.2 CHO Kernels

The kernels in CHO are diverse, substantial, real world applications. They have all been shown to map well to FPGA using different languages, tools and technique making them ideal for benchmarking HLS frameworks targeted at FPGAs. Consequently, the suite is biased towards “embedded computing” where FPGAs have been used for decades for implementing algorithms. We describe briefly, next, the applications in the suite.

dfadd: is an implementation of IEC/IEEE-standard double-precision floating-point addition using

Domain	Application	Description	Original Sources
Arithmetic	dfadd	IEC/IEEE double-precision floating-point addition	SoftFloat [18]
	dfdiv	IEC/IEEE double-precision floating-point division	SoftFloat [18]
	dfmul	IEC/IEEE double-precision floating-point multiplication	SoftFloat [18]
	dfsfn	Double-precision floating-point sine function	SoftFloat [18]
Media	adpcm	Adaptive differential pulse code encoder & decoder	SRTB [19]
	gsm	GSM residual pulse excitation/long term prediction coding	MediaBench [20]
	jpeg	JPEG image decoder	CHStone [4], PVRG [21]
	motion	Motion vector decoding for MPEG-2	AiLab [22]
Cryptography	aes	Implementation of Advanced Encryption Standard	AiLab [22]
	blowfish	Blowfish Encryption Algorithm	MiBench [23]
	sha	Secure Hash Algorithm	MiBench [23]
Miscellaneous	MIPS	Simplified MIPS processors	CHO [4], SoftFloat [18]

TABLE 1
CHO Kernels

- 64-bit integers.
- dfdiv:** is an implementation of IEC/IEEE-standard double-precision floating-point division using 64-bit integers.
- dfmul:** is an implementation of IEC/IEEE-standard double-precision floating-point multiplication using 64-bit integers.
- dfsfn:** implements the sine function using 64-bit integers. It has several functions in common with the previously listed kernels.
- adpcm:** is an implementation of ITU G.722 Adaptive Differential Pulse-Code Modulation (ADPCM) algorithm used in the encoding and decoding audio signals. ADPCM is often used in Voice over IP communications. The kernel includes encoding and decoding functionalities.
- gsm:** is an implementation of Linear predictive coding, a method of encoding good quality speech at a low bit rate. The kernel implements only the lossy sound compression used in GSM (a mobile communication protocol)
- jpeg:** is an implementation of the JPEG still picture compression standards.
- aes:** implements the AES symmetric-key (the same key is used for encrypting and decrypting data) crypto-system. Encryption and decryption modules are implemented.
- blowfish:** is an implementation of the encryption function of blowfish. Blowfish is a symmetric-key block cipher.
- sha:** implements the Secure Hash Algorithm, a cryptographic hash function.

4.3 Functional Verification

We tested CHO to ensure that we correctly implemented each kernel. First, we verified that the kernels can be parsed by any standard-complaint OpenCL 1.0 front-

Device	Type	Driver/Compiler
Intel Core i5	CPU	Intel OpenCL SDK 2013R3
Intel Core i5	CPU	AMD APP SDK 2.9
Intel Xeon E3	CPU	Intel OpenCL SDK 2013R3
Intel Xeon E3	CPU	AMD APP SDK 2.9

TABLE 2
Functional Test Platforms

end. Since OpenCL-based HLS frameworks often rely on LLVM for their front-end processing we used Clang’s [24] OpenCL 1.0 parser for our tests. Clang is a compiler front-end that converts C-based programming languages into LLVM IR. Syntax checking of kernels is important as OpenCL compilers are often varyingly lax in enforcing syntax correctness.

We validated that the kernels are functionally correct and produce the right results by compiling and running all kernels on diverse OpenCL devices. The devices used are shown in Table 2. All the kernels compiled and produced the correct results on execution.

5 CHARACTERIZING CHO

We characterised CHO at the source-level, after conversion into Abstract Syntax Tree (AST) and the IR-level, after conversion to LLVM-IR.

5.1 Source-Level Characterization

We implemented a Clang plugin that walked each kernel’s AST and classified every token. We run the plugin on the AST produced by the Clang compiler. Table 3 is a summary of the source-level characteristics of the each kernel. ‘Dominant Type’ in the table refers the representative data type in each kernel while ‘LoC’ refers to the lines of code

Kernel	Dominant Type	LoC	Functions	Variables		Statements					Operators						
				Scalar	Aggregate	<i>for</i>	<i>if</i>	<i>goto/break</i>	<i>switch</i>	<i>while</i>	Divide	Multiply	Add/Sub	Compare	Shift	Assign	Logic
adpcm	32-bit <i>int</i>	463	21	103	13	10	17	1	0	0	4	36	112	27	29	171	7
aes	32-bit <i>int</i>	667	10	42	25	16	24	31	6	0	14	95	98	41	184	389	209
blowfish	8-bit <i>char</i>	495	4	38	16	6	9	0	0	4	0	9	178	16	123	267	267
dfadd	64-bit <i>int</i>	352	81	2	18	1	43	6	0	0	0	0	29	46	21	104	34
dfdiv	64-bit <i>int</i>	280	19	96	1	1	31	0	0	2	2	4	33	45	30	100	31
dfmul	64-bit <i>int</i>	350	16	74	1	1	28	0	0	0	0	4	25	38	24	90	31
dfsine	64-bit <i>int</i>	611	31	178	1	1	74	6	0	3	2	7	54	91	45	188	65
gsm	32-bit <i>int</i>	310	12	45	7	17	16	0	0	1	0	47	168	57	23	211	11
jpeg	32-bit <i>int</i>	1061	31	204	30	34	58	11	1	11	3	52	135	92	44	342	21
mips	32-bit <i>int</i>	215	1	17	3	2	3	35	3	4	0	2	12	10	22	57	23
motion	32-bit <i>int</i>	436	13	59	15	7	21	0	0	4	2	6	40	28	16	84	7
sha	64-bit <i>char</i>	178	8	42	15	10	2	0	0	4	2	1	42	16	23	90	32

TABLE 3
Source-Level Characteristics

in the source excluding comments and empty lines. The table shows that kernels are non-trivial implementations with hundreds of lines of code, multiple functions plus diverse statements and operators.

Some kernels perform multiplications and divisions which are expensive. On FPGAs, division is the most expensive operator to implement followed by multiplication [25]. Modern FPGAs often have DSP blocks with high-speed embedded multipliers making it possible to multiply without using logic resources. However, there are only a few of these on FPGAs which implies that kernels that perform a lot of multiplications such as **aes** may still have to implement some of the multipliers using logic resources. All divisions must be implemented with logic resources which may use-up logic resources if there are many of them. Generally, these expensive functions can be shared reducing the number required as well as performance.

5.2 IR-Level Characterization

Since LLVM IR instructions are simple enough to directly correspond to hardware operations characterizing at the IR-level provides a more accurate and detailed picture about how kernels may be mapped to FPGAs. For instance, a single loop at the source-level may be simplified and split into multiple loops at the IR-level.

For IR-level characterization, we implemented a Clang/LLVM compiler that mimics the front-end of a generic OpenCL HLS compiler. Our custom compiler first translates OpenCL into LLVM IR and then applies transformations and optimizations [26] that has been shown to improve performance on FPGAs. These optimizations are shown in Table 4.

Table 5 summarizes the IR-level characteristics of CHO kernels. Notice that the number of loops has increased moving from source-level to IR-level as a result of running `simplifycfg`. The compiler splits complex loops into multiple simpler efficient loops. Figure 9 is a breakdown of the IR-level instructions by instruction types. In the diagram, ‘Memory’ refers to load and store type

Optimization & Transformation	Description
<code>inline</code>	Bottom-up inlining of functions into callees
<code>jump-threading</code>	Attempts to find distinct threads of control flow running through a basic block
<code>simplifycfg</code>	Removes dead code and merges basic-blocks
<code>loop-rotate</code>	Simple loop rotation
<code>gvn</code>	Removes fully and partially redundant instructions plus redundant load elimination
<code>instcombine</code>	Combines instructions into fewer and simpler instructions

TABLE 4
LLVM transformations and optimizations

Kernel	Number Basic Blocks	Number Instructions	Number Loops
adpcm	95	2240	19
aes	186	5685	23
blowfish	58	5441	10
dfadd	199	1421	1
dfdiv	115	1212	3
dfmul	90	846	1
dfsine	469	4216	2
gsm	149	1493	18
jpeg	661	10261	109
mips	46	472	3
motion	1743	15121	326
sha	79	1364	30

TABLE 5
IR-Level Characteristics

instructions plus address generating instructions while ‘Control’ refers to instructions, such as branch, that transfer control from one part of a kernel to another. For some kernels e.g. **gsm** the expensive operators, multiplication and division, are still significant.

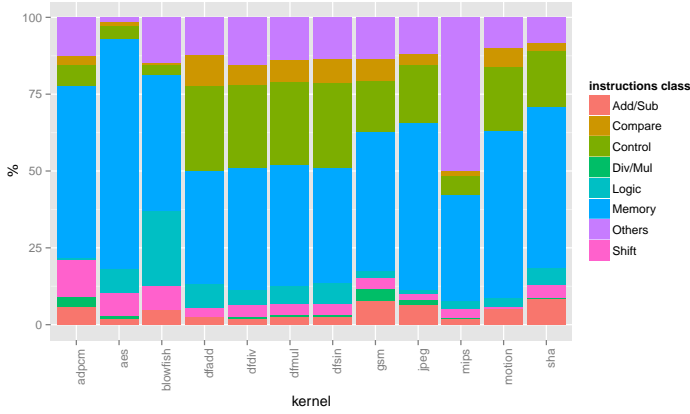


Fig. 9. LLVM IR instruction classes

FPGA	Altera 5SGXMA7H2F35C2
Form Factor	Half-length Half-height PCIe Card
Host Interface	8-lane Gen3 PCIe
Memory	8GB DDR3 SDRAM

TABLE 6
P385-A7 FPGA accelerator card

6 SYNTHESIZING CHO

In this section, we present the results of our attempts at synthesizing kernels in the CHO benchmark suite using AOCL, a state-of-the-art OpenCL HLS framework. The main objective is to determine if we can synthesize unmodified OpenCL kernel for FPGA. We have not modified the kernels so they are the same as those that run on the ‘software’ platforms in Table 2.

6.1 Target Compiler and Platform

We use AOCL 13.1.3.178 and target Nallatech P385-A7 FPGA accelerator card. The P385-A7 (its features are described in table Table 6) is based on Altera’s Stratix V GX-A7 FPGA.

AOCL supports augmenting OpenCL kernels by specifying kernels attributes. The available attributes are listed below.

- `reqd_work_group_size`: Specifies the required work-group size allowing AOCL to allocate the exact amount of FPGA resources needed to manage work-items in a work-group.
- `max_work_group_size`: Defines the maximum number of work-items that can be allocated within a work-group.
- `num_compute_units`: Sets the number of compute units to be instantiated to process the kernel. AOCL distributes the work-groups within a kernel across the specified number of compute units increasing throughput. This may increase global

memory bandwidth contention among compute units.

- `num_simd_work_items`: This is similar to the attribute above but involves only replicating the datapath in a single instruction multiple data (SIMD) fashion. For AOCL to implement a SIMD datapath, the value of `num_simd_work_items` must evenly divides the value specified for `reqd_work_group_size`.
- `max_share_resources`: Limits the number of times an operator e.g. multiplier can be reused without reducing computational throughput.
- `num_share_resources`: Sets the number of times an operator can be reused.
- `max_unroll_loops`: Limits the number of times AOCL can unroll each loop in a kernel. This can be overridden for individual loops using `#pragma unroll`.
- `local_mem_size`: Defines the size of local memory allocated to pointer argument with local specifier. The default is 16 kB.
- `task`: Specifies that the kernel executes in a task parallel fashion.

AOCL applies what is termed resource-driven optimizations to kernels when the `O3` flag is used. This essentially identifies a set of values for `num_compute_units`, `num_simd_work_items`, `max_share_resources`, `num_share_resources` and `max_unroll_loops` that maximizes performance without exceeding resource usage threshold. The default resource usage threshold is 85% but can be changed using the `U` flag.

Table 7 shows the features of Altera Stratix V GX-A7. In the table, Adaptive Logic Modules refers to Altera’s logic blocks. Each block has 8 inputs with a Adaptive Look-Up Table (ALUT) (Altera name for LUTs), 2 adders and 4 registers. Variable Precision Multipliers are the embedded, high-speed, variable precession multipliers. They can perform 9-bit, 18-bit, 27-bit and 36-bit word lengths operations. These multipliers are part of a larger structure called DSP block. Stratix V devices contain dedicated memory blocks (Block RAMs), each 20-Kb in size, called M20K blocks. The number of blocks and the total number of RAM bits available from all M20K blocks are given in the table.

6.2 Sythesis Results

We were unable to synthesize and run every kernel in the CHO for various reasons. For some of the kernels, the LLVM front-end crashed with a stack dump while for others the front-end stopped after an internal compiler error. A few of the Verilog RTL generated couldn’t be synthesized by Quartus, the back-end of AOCL. A summary of our attempts at synthesizing the kernels is presented as Table 8. In the table, ✓ means that a kernel was successfully synthesised.

Adaptive Logic Modules (K)	235
Registers (K)	939
M20K Memory Blocks	2,560
M20K Memory (Mbits)	50
Variable Precision Multipliers (18x18)	512
Variable Precision Multipliers (27x27)	256
DSP Blocks	256
Fabric clock (MHz)	800

TABLE 7
Starix V GX-A7 features

Kernel	Synthesizable	Notes
adpcm	X	Frontend error
aes	X	Frontend crash
blowfish	✓	
dfadd	✓	
dfdiv	X	Frontend crash
dfmul	✓	
dfsine	X	Frontend crash
gsm	✓	
jpeg	X	Frontend error
mips	X	Backend error
motion	X	Frontend crash
sha	X	Backend error

TABLE 8
Summary of Synthesis

The kernels were compiled using the `o3` flag and setting the `task` attribute. Table 9 show the OpenCL attributes used while Table 10 presents the synthesis results. Note that the OpenCL attributes in Table 9 were set automatically by the AOCL.

In Table 10, ‘Loops’ refers to the number of loops at the LLVM-IR level while ‘Pipelined Loops’ refers to the number of those loops that AOCL managed to fully pipeline i.e. loops without Loop-Carried Dependencies (LCD) [17]. AOCL attempts to pipeline [27] (overlapping computations for different loop iterations in time and space) all loops so as to have a fully pipelined hardware structure. Some of the loops may have LCD which requires AOCL to generate extra hardware to account for these memory dependencies reducing throughput. Therefore, rewriting loops (following Altera’s recommendations [28]) to minimize LCD is a good approach to improving performance. Pipelining in the presence of LCD is an active research area and techniques developed in this area could be beneficial to future HLS frameworks.

LCD is particularly a problem for **blowfish** where 3 of its loops each causes, on the average, a 567 cycles pipeline stall (i.e. successive iterations launched every 567 cycles) The cumulative effect of this is a significant degradation of performance as shown in the next section. Note that information about loops are available from the AOCL

Kernel	Attributes Settings
blowfish	max_unroll_loops (2) num_compute_units (2) num_share_resources (1) max_share_resources (8)
dfadd	max_unroll_loops (1) num_compute_units (3) num_share_resources (1) max_share_resources (8)
dfmul	max_unroll_loops (40) num_compute_units (1) num_share_resources (1) max_share_resources (8)
gsm	max_unroll_loops (2) num_compute_units (1) num_share_resources (1) max_share_resources (8)

TABLE 9
OpenCL attributes

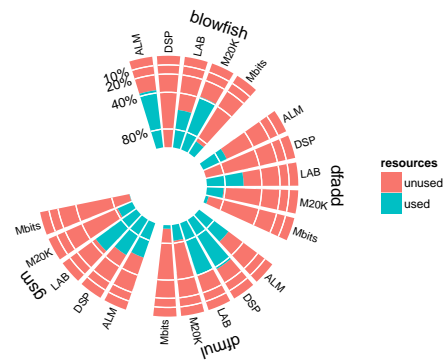


Fig. 10. Percentage resource utilization

during compilation.

Notice that the number of loops in Table 5 and Table 10 don’t exactly correspond. This is because our custom compiler and AOCL may not be using the same set of optimizations and versions. Further, AOCL completely unrolls some of the loops as could be inferred from **dfmul**.

Figure 10 shows the percentage resource usage in terms on ALMs, LABs, DSP blocks, M20K blocks, total number of memory bits (Mbits in the figure). It gives a rough indication of the “complexity” of each implementation on the FPGA.

6.3 Performance

Figure 11 compares the execution time of the successfully synthesized kernels on the FPGA platform and on a CPU platform. The details of the CPU platform is shown as Table 11. Note that OpenCL is not performance portable and that we didn’t manually optimize the kernels for FPGA. Compilation for the CPU was done offline to

Kernel	fmax (MHz)	ALMs	LABs	ALUTs	Registers	M20K	Memory Bits	DSP Blocks	Synthesis Time (Hours:Mins)	Loops	Pipelined Loops
blowfish	200	147 981	19 724	151 907	270 884	1567	7 946 172	0	03:00	11	7
dfadd	246	65 857	9 415	70 991	107 176	468	1 731 152	0	01:15	1	1
dfmul	207	93 927	13 935	96 103	16 752	443	2 050 512	160	01:30	0	0
gsm	263	93 298	12 914	98 391	161 772	407	162 868	116	01:36	11	9

TABLE 10
Synthesis results

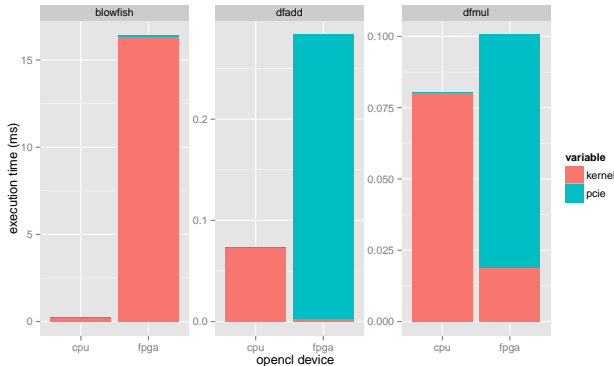


Fig. 11. Comparing performance

CPU	Intel Xeon CPU E31245 @ 3.30GHz
Memory Size	16 GB DDR3
OpenCL SDK	Intel OpenCL SDK 2013R3

TABLE 11
CPU platform features

match AOCL. AOCL presently doesn't support out-of-order execution i.e. the execution of multiple tasks in parallel therefore performance measurements was done with a single task executing on the CPU and the FPGA. We verified that results produced by the FPGA match exactly with the ones from the CPU. We couldn't get **gsm** to run to completion so it doesn't appear in Figure 11.

FPGA outperforms the CPU, without the overhead of PCIe data transfer, except for **blowfish** where the FPGA is slower than the CPU, by an order of magnitude. As mentioned earlier, **blowfish** suffers from having loops that are significantly affected by LCD.

7 CONCLUSION

In this paper, we presented CHO, an OpenCL port of the commonly used CHStone HLS benchmark. CHO aims are to enable HLS framework developers to qualitatively evaluate new ideas as well as enable the benchmarking of the increasing number of OpenCL HLS frameworks.

We characterised CHO at the OpenCL source-level and at the LLVM IR-level. We showed that kernels in the benchmark suite are substantial. IR-level characterization provides more detailed information as the LLVM IR instructions are simple enough to directly correspond to hardware operations.

We attempted to synthesize all the kernels in CHO using AOCL but only managed to produce FPGA bitstreams for 4 out of the 12 kernels. Synthesis failed for a number various reasons ranging from front-end crash to the failure of Quartus, the back-end of AOCL, to synthesize Verilog RTL produced by front-end. We recognize that AOCL is relatively new and like other compilers could improve with time. We showed that when compilation works it is straightforward to compile unmodified OpenCL kernels down to FPGA.

There are a number of possible directions for future work.

- The programs in CHStone, which we ported to OpenCL, are skewed towards "embedded computing" where FPGAs has been used for decades. We are looking at introducing programs from other areas, such databases, that have been shown to have efficient FPGA implementations.
- CHO employs, exclusively, OpenCL's task parallel programming model. However, an OpenCL-based HLS framework may compile the two programming models (slightly) differently as is the case in AOCL. We are looking at introducing kernels into CHO that can be mapped efficiently onto the data parallel programming model of OpenCL.

8 SOFTWARE DOWNLOAD

CHO is publicly available at <http://it302.github.io/cho>.

ACKNOWLEDGEMENT

The authors would like to thank...

REFERENCES

- [1] *The OpenCL Specification Version: 1.0*, Khronos OpenCL Working Group, Std., June 2009.

- [2] K. Morris, "HLS versus OpenCL Xilinx and Altera Square Off on the Future," *Electronic Engineering Journal*, March 2013. [Online]. Available: <http://www.eejournal.com/archives/articles/20130312-highlevel>
- [3] P. Coussy *et al.*, "An Introduction to High-Level Synthesis," *Design Test of Computers, IEEE*, vol. 26, no. 4, pp. 8–17, July 2009.
- [4] Y. Hara *et al.*, "Proposal and quantitative analysis of the chstone benchmark program suite for practical c-based high-level synthesis," *Journal of Information Processing*, vol. 17, pp. 242–254, 2009.
- [5] N. Dutt, "Benchmarks for the 1992 High Level Synthesis Workshop," 1992. [Online]. Available: <http://www.cbl.ncsu.edu:16080/benchmarks/HLSynth92/GUIDELINES>
- [6] P. Panda and N. Dutt, "1995 high level synthesis design repository," in *System Synthesis, 1995., Proceedings of the Eighth International Symposium on*, Sep 1995, pp. 170–174.
- [7] P. Mistry *et al.*, "Valar: A Benchmark Suite to Study the Dynamic Behavior of Heterogeneous Systems," in *Proceedings of the 6th Workshop on General Purpose Processor Using Graphics Processing Units*, ser. GPGPU-6. New York, NY, USA: ACM, 2013, pp. 54–65.
- [8] S. Che *et al.*, "A characterization of the rodinia benchmark suite with comparison to contemporary cmp workloads," in *Workload Characterization (IISWC), 2010 IEEE International Symposium on*, Dec 2010, pp. 1–11.
- [9] A. Danalis *et al.*, "The Scalable Heterogeneous Computing (SHOC) Benchmark Suite," in *Proceedings of the 3rd Workshop on General-Purpose Computation on Graphics Processing Units*, ser. GPGPU '10, New York, NY, USA, 2010, pp. 63–74.
- [10] M. Owaidia *et al.*, "Synthesis of platform architectures from opencl programs," in *Field-Programmable Custom Computing Machines (FCCM), 2011 IEEE 19th Annual International Symposium on*, May 2011, pp. 186–193.
- [11] "Altera SDK for OpenCL: Programming Guide," Dec. 2013.
- [12] D. Koch and J. Torresen, "FPGASort: A high performance sorting architecture exploiting run-time reconfiguration on FPGAs for large problem sorting," in *Proceedings of the 19th ACM SIGDA international symposium on Field programmable gate arrays*. ACM, 2011, pp. 45–54.
- [13] D. B. Thomas, L. Howes, and W. Luk, "A comparison of cpus, gpus, fpgas, and massively parallel processor arrays for random number generation," in *Proceedings of the ACM/SIGDA International Symposium on Field Programmable Gate Arrays*, ser. FPGA '09. ACM, 2009, pp. 63–72.
- [14] H. G. Dietz, *The Aggregate Magic Algorithms*, Aggregate.Org online technical report, University of Kentucky Std. [Online]. Available: <http://aggregate.org/MAGIC/>
- [15] T. Czajkowski *et al.*, "From opencl to high-performance hardware on fpgas," in *Field Programmable Logic and Applications (FPL), 2012 22nd International Conference on*, Aug 2012, pp. 531–534.
- [16] C. Lattner and V. 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.
- [17] V. H. Allan *et al.*, "Software pipelining," *ACM Comput. Surv.*, vol. 27, no. 3, pp. 367–432, Sep. 1995. [Online]. Available: <http://doi.acm.org/10.1145/212094.212131>
- [18] J. R. Hauser, "Softfloat release 2b general documentation," 2007. [Online]. Available: <http://www.jhauser.us/arithmetic/SoftFloat.html>
- [19] SNU Real-Time Benchmarks. [Online]. Available: <http://archi.snu.ac.kr/realtime/benchmark>
- [20] L. Chunho, M. Potkonjak, and W. Mangione-Smith, "Mediabench: a tool for evaluating and synthesizing multimedia and communications systems," in *Microarchitecture, 1997. Proceedings., Thirtieth Annual IEEE/ACM International Symposium on*, Dec 1997, pp. 330–335.
- [21] A. C. Hung, "PVRG-JPEG CODEC 1.1," Stanford University, Tech. Rep., 1993.
- [22] "Ailab." [Online]. Available: <http://www-ailab.elcom.nitech.ac.jp>
- [23] M. Guthaus *et al.*, "Mibench: A free, commercially representative embedded benchmark suite," in *Workload Characterization, 2001. WWC-4. 2001 IEEE International Workshop on*, Dec 2001, pp. 3–14.
- [24] "Clang." [Online]. Available: <http://clang.llvm.org>
- [25] J. Liu, M. Chang, and C.-K. Cheng, "An iterative division algorithm for fpgas," in *Proceedings of the 2006 ACM/SIGDA 14th International Symposium on Field Programmable Gate Arrays*, ser. FPGA '06. New York, NY, USA: ACM, 2006, pp. 83–89.
- [26] Q. Huang *et al.*, "The effect of compiler optimizations on high-level synthesis for fpgas," in *Field-Programmable Custom Computing Machines (FCCM), 2013 IEEE 21st Annual International Symposium on*, April 2013, pp. 89–96.
- [27] D. D. Gajski *et al.*, *Embedded System Design: Modeling, Synthesis and Verification*. Springer, 2009.
- [28] "Altera SDK for OpenCL: Optimization Guide," Dec. 2013.