Evaluating Reconfigurable Dataflow Computing Using the Himeno Benchmark

Yukinori Sato (JAIST, JST CREST) Yasushi Inoguchi (JAIST) Wayne Luk (Imperial College London) Tadao Nakamura (Keio University)

2012 International Conference on ReConFigurable Computing and FPGAs 5th December, 2012







•JAIST •NAGOYA •OSAKA

AIST

Japan Advanced Institute of Science and Technology

Introduction

Acceleration using FPGA

- An approach to build custom system appropriate for each application
- Realize custom system within a certain power, or cost budget
- Customization and optimization must be performed

Customization and optimization strategies must be made based on the characteristics of platform and the type of applications

Characteristics of platforms	Hardware Reconfiguration capability, Customizability, Device itself	<mark>Software</mark> System software, HW compiler, API, run-time, library
Types of applications	Computation-bound apps The clock rate or parallelism should be increased	Memory-bound apps Optimizations for memory is the most significant factor

AIST

In this paper

We focus on FPGA accelerator provided by Maxeler, and evaluate its performance for a memory-intensive application

- We present several optimization techniques
 - Targeting for Maxeler's reconfigurable dataflow engines
 - In terms of memory locality and memory bandwidth

We demonstrate all key optimizations can be done by high-level programming for the Maxeler platform

- For evaluation, we use the Himeno Benchmark
 - Widely known as a memory-bound FP program
- We compare it with the implementation on GPUs

3

The Himeno Benchmark

- Developed by Dr. Ryutaro Himeno, RIKEN, Japan
- A highly memory intensive application kernel
 - Becoming popular in the HPC community to evaluate the worst-case performance for memory bandwidth intensive codes
 - The kernel is a linear solver for 3D pressure Poisson equation which appears in a Navier-Stokes solver:

$$\frac{\partial^2 p}{\partial x^2} + \frac{\partial^2 p}{\partial y^2} + \frac{\partial^2 p}{\partial z^2} + \alpha \frac{\partial^2 p}{\partial xy} + \beta \frac{\partial^2 p}{\partial xz} + \gamma \frac{\partial^2 p}{\partial yz} = \rho$$

- This Poisson equation is solved using the Jacobi iterative method
- The performance of this is measured in FLOPS (FLoating-point Operations Per Second)

4

Source code of the Himeno benchmark



The body of this computational kernel involves 34 FP ops So, FLOPS can be measured by dividing it by the total execution time

AIST

Reconfigurable dataflow engine

- Dataflow graph is directly mapped into FPGA
- Computations are performed by forwarding intermediate results directly from one functional unit to the next units
- Memory transfer operations can be dramatically reduced compared with instruction-based processor



An overview of dataflow processing: Its datapath becomes a deeply pipelined structure

Using inherent locality of dataflow computation, we attempt to avoid memory bottleneck

AIST

An outline of Maxeler's dataflow engine

A high-level HW synthesis platform

- Kernel code is generated by MaxCompiler
 - Use Java as a meta-programing language
 - The structure of dataflow graph is described
 - Java extensions are used for hardware descriptions
 - MaxCompiler generates FPGA bitstream
 - Data exchange between a host CPU and a dataflow engine is performed using a run-time library API (called MaxCompilerRT)
- The bitstream file (.max) is loaded from the CPU to the engine, and the FPGA is configured to the design

Using a high-level description language for Maxeler platform, we optimize design focusing on memory locality and bandwidth

O. Pell and O. Mencer, "Surviving the end of frequency scaling with reconfigurable dataflow computing," ACM SIGARCH Comput. Archit. News, vol. 39, no. 4, pp. 60–65, Dec. 2011.



Optimization strategies for dataflow engine

We present the following 4 optimization strategies

Optimization (A):

Directly forwarding data

Making use of <u>fine-grained parallelism</u> in dataflow graph, <u>temporal parallelism</u> by pipelining dataflow paths

Optimization (B):

Organizing highly regular data stream

The regularity of the flow of data is strongly emphasized so that <u>deeply</u> <u>pipeline processing can be performed</u> without fatal pipeline stalls

Optimization (C):

Utilizing available HW resources

By <u>unrolling independent loops</u>, we build <u>multiple parallel pipes</u> in a kernel

Optimization (D):

Communication btw CPU and FPGA

When <u>partitioning the code regions for</u> <u>an FPGA accelerator</u>, the amount of communication must be investigated carefully

Based on these 4 strategies together with domain knowledge and experience for apps, we optimize the kernel

Japan Advanced Institute of Science and Technology

8

AIST

Optimizing the Himeno benchmark (1)

- Perform <u>Optimization (A)</u> to improve locality of memory access
 - We cascade multiple sequential points of the stencil computation
 - This enables to keep more intermediate results within pipelines and mitigates heavy memory traffic of a stencil program



Stencil window (7-point)

Cascading the next point

An overview of cascading:

Since sequential points of a stencil computation often access the same data, we can reuse points appeared previously

Japan Advanced Institute of Science and Technology

9

Optimizing the Himeno benchmark (2)

- Perform Optimization (B) by converting data stream into 1-D
 - To create a window for stencil computation, we need to access the adjacent values in a 3D array structure
 - We use stream.offset (Maxeler API) to feed data into a kernel
 - MaxCompiler transforms the stream.offset into an FIFO buffer on FPGA



Converting 3-D data accesses into 1-D stream:

By specifying a relative distance from the current stream data, we can obtain the past or the future data.

Japan Advanced Institute of Science and Technology

TAIST

Optimizing the Himeno benchmark (3)

- Perform <u>Optimization (C)</u> by configuring the # of pipes
 - We adjust the # of pipes to maximize parallelism and resource utilization of an FPGA chip
 - Here, we refer the body of a custom dataflow pipeline as a **pipe**
 - By unrolling loop iterations based on degree of cascading of stencil computations, we realize the parameterized kernel



Japan Advanced Institute of Science and Technology

Optimizing the Himeno benchmark (4)

• Perform <u>Optimization (D)</u> by building 3 scenarios for the stream communications



PCIe-** design

- Data stream via PCIe bus
- Simple
- But data transfer via PCIe bus will be bottleneck





nn-itr-** design

- Data stream via internal buffer
- Consume BRAM equal to the size of array 'p'

DRAM-** design

- Data stream via on-board DRAM
- Require memory adr gen

JAIST

Japan Advanced Institute of Science and Technology

Performance evaluation

Methodology

Maxeler MaxWorkstation

- Intel Core i7-870 2.93GHz CPU, CentOS
- The MAX3 acceleration card
 - Virtex-6 SX475T FPGA, 24GB DDR3 memory
 - PCI express gen2 x8 interface
- MaxCompiler version 2012.1
- Generate host code using gcc.4.1.2 with '-O3'

Experimental conditions

- Set FPGA operating frequency as 100MHz (system default)
- Constant array values are generated using ternary-if logic







Evaluation results

Verification of our designs

- Compare the output values of p from the FPGA with those of the original CPU implementation across all of design
- As a result, we confirmed that all of these two are completely the same → [accurate implementation]
- Also, we could implement our design smoothly

The productivity of Maxeler platform with support for IEEE floating point format enables *accurate* and *smooth* implementations



Stream communication via PCIe



Performance result		
# of pipes	score [GFLOPS]	
1	2.70	
2	4.96	
4	8.29	
8	8.33	

PCIe-** desing

- Communication via PCIe bus
- The data stream is driven by the triply-nested i,j,k loop
- Simple implementation
- But data transfer via PCIe bus will be bottleneck

- From the result, we find that the perf. is not increased even if we scale the # of pipes
- The data transfer via the PCIe bus becomes the bottleneck



Data stream via an internal buffer

In order to avoid the data transfer overhead, we attempt to feed output data directly to the input



• The data stream is driven by the outermost loop

- For keeping the values from the previous loop iteration, we use stream.offset for inserting a special buffer
- Once the host CPU transfers the initial data stream, data stream from the previous iteration is fed into the pipes directly



The results obtained by data stream via an internal buffer



Evaluation Results

- This design can linearly increase its performance when we scale the # of pipes
- The data transfer bandwidth bottleneck can be avoided
- We successfully build and run the design with 48 pipes
- The 48-pipe design with 110MHz achieves 155.1 GFLOPS

Due to the size of BRAM, configurations with larger data set cannot be implemented

TAIST

Data stream via on-board DRAM

In order to apply larger data sets, we attempt to use on-board DRAM devices



DRAM-** design

- The data stream is driven by the triply-nested i,j,k loop similar to the 'PCle-**' design
- This design can make use of higher on-board DRAM memory bandwidth compared with the PCIe design
- Iterations are processed without any data transfer to the CPU memory once the initial input data are loaded



The results obtained by data stream via on-board DRAM



Evaluation Results This design can accept larger data sets (up to L: 129.3MB) and increase its performance according to the # of pipes We cannot build the 48 pipe designs due to HW overhead for memory address generators We can see start-up overhead of the kernel pipelinein smaller data set

• 32-pipe design with L data set can achieve 97.6 GFLOPS

TAIST

Comparison with the state-of-the-art GPU



[Phillips 2010] E. Phillips and M. Fatica, "Implementing the himeno benchmark with CUDA on GPU clusters," in *IEEE Int'l Symp. on Parallel Distributed Processing*, 2010, pp. 1–10.

AIST

Japan Advanced Institute of Science and Technology

Conclusions

We describe a memory bound application program on an FPGA accelerator based on a reconfigurable dataflow computing platform

- We have presented several optimization techniques for designs targeting reconfigurable dataflow engines
 - especially focus on memory locality and memory bandwidth
 - Also present how such techniques can be used in optimizing the Himeno Benchmark that requires large memory bandwidth
- From the results of evaluation
 - We have confirmed that specialized dataflow pipeline designs contribute to improving the actual performance of applications
 - Found that our implementation can outperform the recent GPU • implementations in their achieved GFLOPS
 - Demonstrated that we can achieve competitive performance gain without low-level HDL coding



AIS

Future work & Acknowledgment

Current and future work

- Extending the set of optimization strategies and benchmarks for our approach
- Exploring methods for automating the application of such optimization strategies
- Studying how our approach can be improved to support power and energy optimization

<u>Acknowkedgment</u>

This work was supported in part by the JSPS "Institutional Program for Young Researcher Overseas Visits", by UK EPSRC, by the European Union Seventh Framework Programme under Grant agreement number 248976, 257906 and 287804, by the HiPEAC NoE, by the Maxeler University Program, and by Xilinx.

