

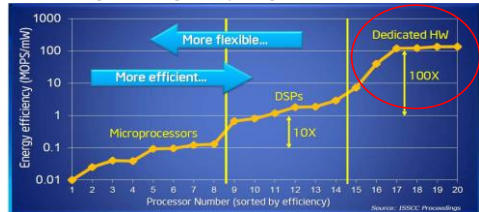
## Designing Application-Specific Heterogeneous Architectures from Performance Models

Minh Thanh CONG - Irisa / Université de Rennes 1  
François CHAROT - Inria  
October 03, 2019



### Processor Design Trends

- Potential of application-specific accelerators
- 10x to 1000x performance per watt gain improvement



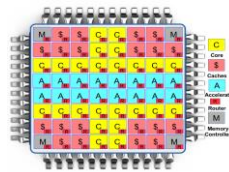
9/25/2019 Source: Bob Broderson, Berkeley Wireless group  
Thanh CONG - University of Rennes 1

### Designing Heterogeneous Architectures

- How can we design such architectures?**
- Simulation: a way to explore the design space
  - Reusing knowledge and design tools of the multi-/many-core domain

- Existing works**
- Standalone accelerator: Aladdin [Shao, 2014]
  - System-level simulation: gem5+Aladdin [Shao, 2016], PARADE [Cong, 2015]
  - Cycle-accurate simulators are still quite slow

- What can we do?**
- Speeding up the cycle-accurate simulator
  - Building performance models of accelerators



### Harnessing FPGAs for Simulation not Prototyping

Implement (parts of) the simulator on an FPGA

- Parallelism
- FPGA attachment technology
- FPGA is more and more powerful



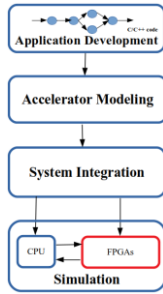
**Performance model**

- Functionally equivalent and logically isomorphic
- Allowing abstractions which simplify model development, enabling modularity

Complexity ↓

### Our contribution

- A methodology for generating performance models of accelerators on FPGAs
- A cycle-accurate simulator targeting CPU-FPGA platforms with the goal of speeding up the simulation



9/25/2019

Thanh CONG - University of Rennes 1

5

### Outline

1. Accelerator modeling flow
  - Flow explanation by an example
2. System Integration
  - Simulation framework
  - Experimental study
3. Conclusion & Future work

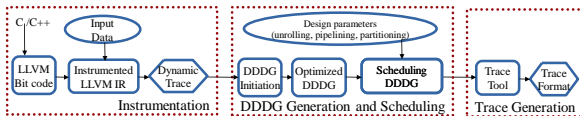
9/25/2019

Thanh CONG - University of Rennes 1

6

### Accelerator modeling flow

- Trace-based simulation
- Constructing a dynamic data dependence graph (DDDG)
- Optimizing and scheduling the graph
- Transforming the scheduled graph into compact instructions



9/25/2019

Thanh CONG - University of Rennes 1

7

### Flow explanation by an example

```

for(i=0; i<N; ++i)
  c[i] = a[i] + b[i];
          
```

A. C Code

```

0. r0=0 //i = 0
1. r4=load(r0 + r1) //load a[i]
2. r5=load(r0 + r2) //load b[i]
3. r6=r4 + r5
4. store(r0 + r3, r6) //store c[i]
5. r0=r0 + 1 //++i
6. r4=load(r0 + r1) //load a[i]
7. r5=load(r0 + r2) //load b[i]
8. r6=r4 + r5
9. store(r0 + r3, r6) //store c[i]
10. r0 = r0 + 1 //++i
...
          
```

B. IR Trace:

C. Un-optimized DDDG

Parameters	Configuration
Loop unrolling factor	2
Loop pipeline	Off
Memory partitioning	Cyclic, factor 2

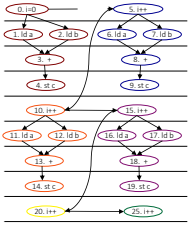
9/25/2019

Thanh CONG - University of Rennes 1

8

### Accelerator Model Targeting FPGAs

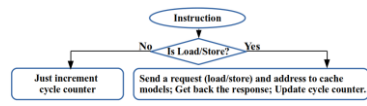
- A factor of 2 loop iteration parallelism, partitioning factor 2, and without loop pipelining



D. Scheduled DDDG  
9/25/2019

Nodes	Instruction	Description
	7 bits: Op_code	
	1 bit: Is new cycle	
Address	32 bits	op_code = load/store

#### E. Accelerator instructions



#### F. Timing model of accelerator

Thank CONG - University of Rennes 1

9

### Outline

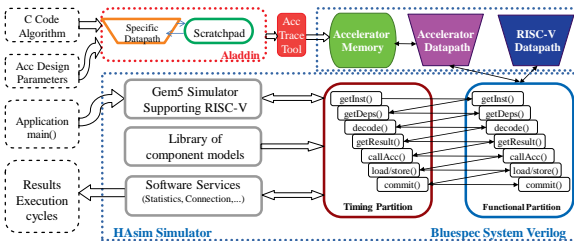
1. Accelerator modeling
  - Flow explanation by an example
2. System Integration
  - Simulation framework
  - Experimental study
3. Conclusion & Future work

9/25/2019

Thank CONG - University of Rennes 1

10

### Simulation Framework

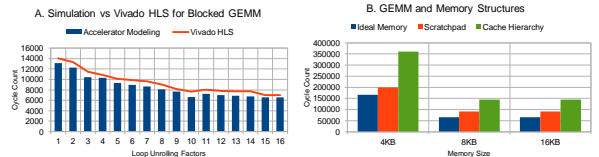


9/25/2019

Thank CONG - University of Rennes 1

11

### Framework Validation Using Blocked GEMM



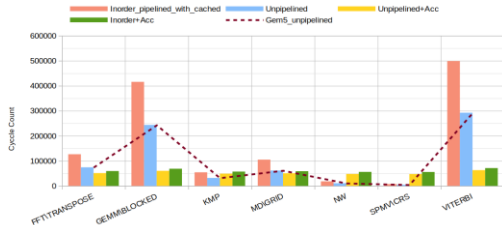
- The latency decreases when the loop unrolling factor increases.
- The performance is very close to Vivado HLS.
- The number of cycles increases as the complexity of memory grows.

9/25/2019

Thank CONG - University of Rennes 1

12

## Architecture Exploration



- Using MachSuite benchmarks
- Same level of precision as the gem5 simulator is observed.

9/25/2019

Thanh CONG - University of Rennes 1

13

## Outline

1. Accelerator modeling
  - Flow explanation by an example
2. System Integration
  - Simulation framework
  - Experimental study
3. Conclusion & Future work

9/25/2019

Thanh CONG - University of Rennes 1

14

## Conclusion

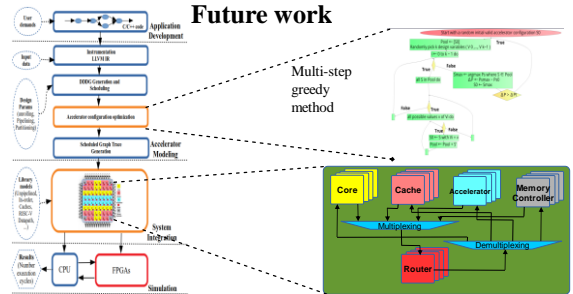
- Generation of accelerator models targeting FPGAs
  - Exploited in HASim-based simulation framework
- Simulation framework allowing architecture exploration
  - Tested with single-core and accelerator models
  - Compared with other tools

9/25/2019

Thanh CONG - University of Rennes 1

15

## Future work



9/25/2019

Thanh CONG - University of Rennes 1

16

## References

- [Shao, 2014] Y. S. Shao, B. Reagen, G. Wei, and D. Brooks, "Aladdin: A pre-rtl, power-performance accelerator simulator enabling large design space exploration of customized architectures," in 2014 ACM/IEEE 41st International Symposium on Computer Architecture (ISCA), pp. 97–108, June 2014.
- [Shao, 2016] Y. S. Shao, S. L. Xi, V. Srinivasan, G. Y. Wei, and D. Brooks, "Codesigning accelerators and soc interfaces using gem5-aladdin," MICRO, 2016.
- [Cong, 2015] J. Cong, Z. Fung, M. Gill, and G. Reinman, PARADE: A cycle-accurate full-system simulation platform for accelerator-rich architectural design and exploration. In Proceedings of the IEEE/ACM International Conference on Computer-Aided Design, ICCAD '15, 2015.
- [Chen, 2016] Yu-Ting Chen, Jason Cong, Zhenman Fung, and Peipei Zhou, 2016. ARAPrototyper: Enabling Rapid Prototyping and Evaluation for Accelerator Rich Architecture (Abstract Only). In Proceedings of the 2016 ACM/SGEDA International Symposium on Field-Programmable Gate Arrays (FPGA '16). ACM, New York, NY, USA, 201-201.
- [Pollauer, 2011] M. M. I. Pollauer, HAsim : cycle-accurate multicore performance models on FPGAs. Thesis, Massachusetts Institute of Technology, 2011.
- [Intel] <https://www.altera.com/solutions/acceleration-hub/overview.html>
- [Xilinx] <https://www.xilinx.com/products/design-tools/vado/integration/cad-design.html>

9/25/2019

Thanh CONG - University of Rennes 1

17

## Questions



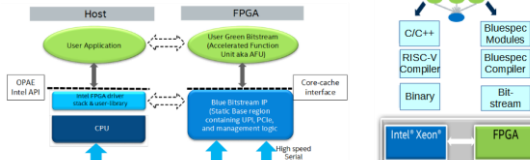
9/25/2019

Thanh CONG - University of Rennes 1

18

## Intel Xeon-FPGA platform

- The Intel® Xeon®-FPGA tightly-coupled FPGA platform
- Minimizing communications bottlenecks



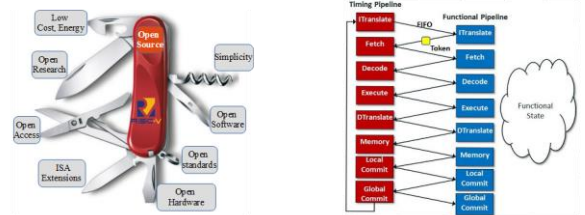
Source: <https://opae.github.io>

9/25/2019

Thanh CONG - University of Rennes 1

19

## RISC-V processor models



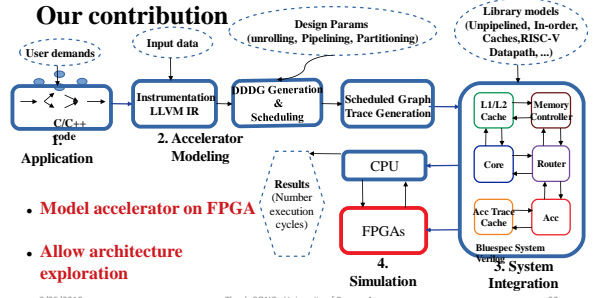
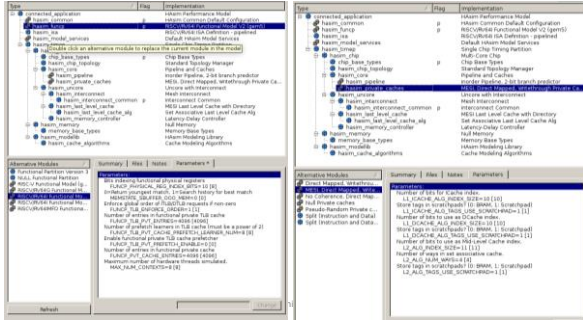
- Support RISC-V 32, 64 bits Integer Instruction Set Architecture

9/25/2019

Thanh CONG - University of Rennes 1

- Build two timing models: unpipelined and in-order pipelined

20



- Model accelerator on FPGA
- Allow architecture exploration

9/25/2019

Thanh CONG - University of Rennes 1

22