Embedded Intelligent System and Novel Computer Architecture

Lecture 06 – FPGA (Field Programmable Gate Array)

and CGRA (Coarse-Grained Reconfigurable Architectures)

Pengju Ren Institute of Artificial Intelligence and Robotics Xi'an Jiaotong University

http://gr.xjtu.edu.cn/web/pengjuren

Three Important Trends

Requirement: Big Data Advances in ML, Data Analytics (Challenges: Data-driven discovery, Search, Analyze data in real time)

Limitations: Moore' Law Dennard Scaling (Power\Memory\Utilization Wall)

Challenges: Algorithms change rapidly High NRE Costs with ASICs Increasing compute & memory requirements

Specialization: Domain Specific Architecture Performance / Watt is key

Flexible Hardware

Flexibility: Instructions



[1] Mark Horowitz, Computing's Energy Problem (and what we can do about it), ISSCC 2014

[2] Hameed et al, Understanding Sources of Inefficiency in General-purpose Chips, ISCA 2010

[3] Leng et al, GPUWattch: Enabling Energy Optimizations in GPGPUs, ISCA 2013

Flexibility: Reconfigurable Hardware



- FPGAs, CGRAs are pre-fabricated silicon devices that can be reprogrammed using a stream of configuration bits.
 - Sea of reconfigurable elements
 - Programmable interconnect
 - Statically programmed
 - No instruction overheads

Fine-grained v.s Coarse-grained RA



What is an FPGA ?

FPGA = Field Programmable Gate Array



Programmable logic, Clock, Interconnections and Routing

Programmable in System (ISP)

Dedicated Blocks: Memory\Clock Control\DSP blocks\Embedded Processor\I/O blocks

What is an FPGA ? – Configurable Logic Block



CLBs contains: LUTs for creating arbitrary combinatorial logic functions flip-flops for clocked storage elements,

multiplexers to route the logic within the block and to and from external resources The muxes also allow polarity selection and reset and clear input selection.

Look-Up Table——the Key to re-programmability



Reading(left) and writing(left) of 4-input LUT

- An n-input LUT can be used to implement an arbitrary Boolean-valued function with up to n Boolean arguments
- LUTs can also be used as memory elements, small FIFO, Shift Registers

What is an FPGA ? – Distributed Memories



The FPGA fabric includes embedded memory elements that can be used as random-access memory (RAM), read-only memory (ROM), or shift registers. A single BRAM block can hold a few kilobytes of data (e.g., 4 KiB), a few hundred BRAMs can be accessed in parallel.
 BRAMs can be used for clock domain crossing and bus width conversion in an elegant way.

What is an FPGA ? - DSP



Xilinx DSP48E slice has three input ports (which are 25 bits, 18 bits, 48 bits wide) and provides a 25x18-bit multiplier in combination with a pipelined second stage that can be programmed as 48-bit substractor or adder with optional accumulation feedback.
 DSP units can be used in a variety of modes, and perform operations such as multiply, multiply-and-accumulate, multiply-and-add/subtract, three input addition, wide bus multiplexing, barrel shifting, etc.

What is an FPGA ? – Configurable Connections



Distributed Connection

The programmable routing in an FPGA provides connections among logic blocks and I/O blocks to complete a user-designed circuit. It consists of wires and programmable switches

What is an FPGA ? – Configurable I/O



The I/O block (IOB) is used to drive signals to the pins of the CPLD device at the appropriate voltage levels with the appropriate current. Two main classes of I/O standards being single-ended (used, e.g., in PCI) and for higher performance differential (used, e.g., in PCI Express, SATA, 10G Ethernet, etc.)

FPGA v.s GPU



- Limited Resources
- Develop using HDL(Hardware Description Language)

FPGA v.s CPU and GPU

A processor is programmed with instructions (CPU and GPU)

FPGA contains configurable blocks with logics and configurable connection lines between these blocks, it is programmed with a circuit description. Pros and Cons:

- Low-level hardware Control and Data movement Operations (Low Programming Efficiency)
- Flexibility comes at the cost of large compile (place/route) and debug times
- FPGA Engineers are hard to hire



Presented (a na se					
Section 1	and an and a state	10-140-01-16-			a free to the track free	
Committy of the owner of the	10 M 11 10	10 10 11 10	N		19.48.10 Sc 10	18 AC 81 11 5
and the second s	10.000	18 30 30 -10	ni ut m at	- 82 - 141 - 80 - 141 - 9	e 1/3 ee ui ju	H
A CONTRACTOR OF A	24					Deft.
Consideration of the						
Committee 10						
Carponent II						
Constitution of the						
Constitution of the						
Print Marcos N						
(International Contention of the						
Contractive and						
Contraction of the						
Caratanana IC						
Province N						
Contraction of the						
Concession, etc.	1					
Constanting of the						
Contraction of the						
And I wanted in the						
In successive statements where the	Training and the second			and a second sec		
Ter 10	try. Manufacture and					and the second second
1	13					
2 313	11 Hall			1 i		-

FPGA v.s CPU and GPU

Feature	Analysis	Winner
Floating-point Processing	The total floating-point operations per second of the best GPUs are higher than the FPGAs' with the maximum DSP capabilities.	GPU
Timing Latency	Algorithms implemented into FPGA provide deterministic timing, with latencies one order of magnitude less than GPUs.	FPGA
Processing / Watt	Measuring GFLOPS per watt, FPGAs are 3-4 times better Although still far away, latest GPU products are dramatically improving the power burning.	FPGA
Interfaces	GPUs interface via PCIe, while FPGA flexibility allows connection to any other device via - almost- any physical standard or custom interface.	FPGA
Backward Compatibility	Software developed for older GPUs will work in the new devices. FPGA HDL can be moved to newer platforms, but with some reworking.	GPU
Flexibility	FPGA lacks flexibility to modify the hardware implementation of the synthesized code, being a no-problem issue for GPUs developers.	GPU
Size	FPGA's lower power consumption requires less thermal dissipation countermeasures, implementing the solution in smaller dimensions.	FPGA
Development	Many algorithms are designed directly for GPUs, and FPGA developers are difficult and expensive to hire.	GPU

FPGA Design Flow

Refining your design



Anatomy of a basic CGRA

Hardware Building blocks: Compute, Memory, Interconnect

Compute: ALUs of varying capability

Memory: Programmer-managed scratchpads, caches

□ Interconnect: Statically programmed paths vs. dynamically routed data

Hardware Organization: Topology

Data path hierarchy: ALUs vs. clusters of ALUs

Communication granularity: bit-level vs. word-level

□ Interconnect topology: Mesh, Torus, ...

Software: Programming Model

□ Software abstraction: Threads, VLIW, spatially configurable ALUs, ...

Compiler technology to map high-level applications to CGRAs

CGRA v.s FPGA

Fine-grained reconfigurability introduces overheads Much higher area, delay and power vs. standard cell ASIC Introduce coarse-grained building blocks, much less interconnect

Programmability

Fine grained architecture leads to long place and route times (> 2 hours) Not a good computing substrate target for a compiler

CGRA architecture

New reconfigurable architectures and new compiler technology

Top-Down Design of CGRA

Observation: We can abstract key software constructs that are amenable to hardware acceleration 2

- Nested data and pipeline parallelism
- Data locality
- Parallel Patterns: Software abstractions that capture parallelism and locality
 - Loops with special properties
 - Expressive over wide range of domains (ML, SQL, Graph analytics, etc)
 - Enables building optimizing compilers with aggressive compiler optimizations
- Design a CGRA to accelerate parallel patterns

CGRA: Basic Structure and Configuration

	┣—		Data memory (Scratchpad)			4 D	4. D	
Host Controller 1	.			V _{ddH} V _{ddL}	A + B	A >> B	A > B	(A+B)>>C
		Dry	$PE \leftrightarrow PE \leftrightarrow PE \leftrightarrow PE$	ALU MEM REGS	A - B	$\mathbf{A} + (\mathbf{B} \gg \mathbf{C})$	$\mathbf{A} == \mathbf{B}$	(A+B)< <c< th=""></c<>
					A & B	A	A < B	(A-B)>>C
		emo			A B	A + (B< <c)< td=""><td>A >= B</td><td>(A+B)<<c< td=""></c<></td></c)<>	A >= B	(A+B)< <c< td=""></c<>
		Context m			A^B	A - (B< <c)< td=""><td>A <= B</td><td>A×B_H</td></c)<>	A <= B	A×B_H
					A ~^ B	A-B	A != B	Clip (A , -B , B)
					~A	(A>>C)-B	A - (B>>C)	Clip (A , 0 , B)
					A << B	A×B_L	(A< <c)-b< th=""><th>C?A:B</th></c)-b<>	C?A:B
L				001	Functions of ALU			

Reconfigurable hardware

- Implement hardware structures dynamically
- Objective: high performance but low power
- Key Tech: spatial parallel computing like 2D-PE array
- Connection is programmable → routing

CGRA: Basic Structure and Configuration





- Compilation tool flow is critical to CGRAs
- Mapping & Scheduling: From DFG to Spatial/Temporal deployment



Operator Mapping
 Memory Mapping
 Interconnection configuration











CGRA: orchestrate dataflow and hardware



- Exploiting operator level parallelism
- Finding better OP-PE binding according to CGRA's arch features

Computation models of CGRAs



Single/Multiple Configuration Single/Multiple Data (SCSD, SCMD, MCMD)

Configuration-1 through configuration-3 are independent and asynchronous; rectangles with different colors represent different configurations, and blank ones represent idle

可重构计算处理器

- 可重构计算芯片以运算单元阵列为核心部件,由配置流和数据流共同驱动(不使用指令),从而实现软硬件双编程。
- 应用任务中的每个运算都可以在运算阵列中找到对应的单元,单元间的 互连与任务——对应,由此获得比拟专用电路的能量效率。



Next Lecture : DNN Acceleractor & HiPU Arch (Given by Prof. Wenzhe Zhao)