ECE 6775 High-Level Digital Design Automation Fall 2024

Hardware Specialization



Cornell University



Announcements

- First reading assignment
 - A. Boutros and V. Betz, "<u>FPGA Architecture</u>: <u>Principles and Progression</u>", IEEE CAS-M 2021
 - Complete reading before Thursday 9/5
- Lab 1 and an HLS tool setup guide will be released soon (by Monday)

Recap: Our Interpretation of E-D-A

Exponential in complexity (or Extreme scale)

Diverse increasing system heterogeneity

Algorithmic intrinsically computational



Significance of EDA: Another Proof

of Customs over Time



Ruchir Puri, High Performance Microprocessor Design, and Automation: Challenges and Opportunities, TAU'2013 keynote.

Agenda

- Motivation for hardware specialization
 - Key driving forces from applications and technology
 - Main sources of inefficiency in general-purpose computing
- A taxonomy of common specialization techniques
- Introduction to fixed-point types

A Golden Age of Hardware Specialization

Higher demand on efficient compute acceleration, esp. for machine learning (ML) workloads



Lower barrier with open-source hardware & accelerators in cloud coming of age



Traditional sw (CPU) server plane

Rising Computational Demands of Emerging Applications

- Deep neural networks (DNNs) require enormous amount of compute
 - Consider ResNet50, a 70-layer model that performs 7.7 billion operations to classify an image (a relatively small model by today's standards)



On Crash Course with the End of "Cheap" Technology Scaling



Hennessy and Patterson, Turing Lecture 2018, overlaid over "42 Years of Processors Data" https://www.karlrupp.net/2018/02/42-years-of-microprocessor-trend-data/; "First Wave" added by Les Wilson, Frank Schirrmeister Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2017 by K. Rupp

Dennard Scaling in a Nutshell

- Classical Dennard scaling
 - Frequency increases at constant power profiles
 - Performance improves "for free"!

Dennard scaling

Transistor (trans.) #	S ²
<u>Capacitance / trans.</u>	1/S
Voltage (V _{dd})	1/S
<u>F</u> requency	S
Total power	1

Note: Dynamic power $\propto CV^2F$

End of Dennard Scaling and its Implications

- Power limited scaling
 - V_{th} scaling halted due to exponentially increasing leakage power
 - V_{DD} scaling nearly stopped as well to maintain performance

Leakage limited sc	aling			
Transistor (trans.) #	S ²			
<u>Capacitance / trans.</u>	1/S			
<u>V</u> oltage (V _{dd})	~1			
<u>F</u> requency	~1			
Total power	S			
Note: Duramia nouver of CV/2E				

Note: Dynamic power $\propto CV^2$ F

- Implication: "Dark silicon"?
 - Power limits restrict how much of the chip can be activated simultaneously
 - No longer 100% without more power

Trade-off Between Flexibility and Efficiency



CPU Core Architecture

- Core = complex control + limited # of compute units + large caches
 - Scalar & vector instructions
 - Backward compatible ISA
 - Complex control logic: decoding, hazard detection, exception handling, etc.



- Mainly optimized to reduce latency of running serial code
 - Shallow pipelines (< 30 stages)
 - Superscalar, OOO, speculative execution, branch prediction, prefetching, etc.
 - Low throughput, even with multithreading

Poll & Discussion



Shallow vs. Deep Pipelining in context of CPU design



http://pollev.com/ece6775 Sign in or register using your Cornell email

Multi-Core Architecture

Control	ALU	ALU	Control	ALU	ALU	Control	ALU	ALU ALU	Control	ALU	ALU
Control	ALU	ALU	Control	ALU	ALU	Control	ALU	ALU	Control	ALU	ALU
Private L1/L2 Cache Private L1/L2 Cache				Private L1/L2 Cache Private L1/L2			1/L2 Ca	che			
Shared Last-Level Cache (LLC)											

With four cores, should we expect a 4x speedup on <u>an arbitrary application</u>?

Graphics Processing Unit (GPU)

- GPU has thousands of cores to run many threads in parallel
 - Cores are simpler (compared to CPU)
 - No support of superscalar, OOO, speculative execution, etc.
 - ISA not backward compatible
 - Amortize overhead with SIMD + single instruction multiple threads (SIMT)
- Optimized to increase throughput of running <u>data-parallel applications</u>
 - Initially targeting graphics code
 - Latency tolerant with many concurrent threads





It's Not Just About Performance: Computing's Energy Problem

	Reading two 32b words from	Energy	
	DRAM	1.3 nJ	
	Large SRAM (256MB)	58 pJ	
65,000x	Small SRAM (8KB)	5 pJ	
	Moving two 32b words by	Energy	
	40mm (across a 400mm2 chip)	77 pJ	
	1mm (local communication)	1.9 pJ	250x
	Arithmetic on two 32b words	Energy	
	FMA (float fused multiply-add)	1.2 pJ	
	IADD (integer add)	0.02 pJ	Ϋ́

Data from [1], based on a 14nm process

Data supply far outweighs arithmetic operations in energy cost

[1] William Dally and Uzi Vishkin, On the Model of Computation, CACM'2022.

Rough Energy Breakdown for an Instruction



Principles for Improving Energy Efficiency

Do less work!

 Amortize overhead in control and data supply across multiple instructions

Amortizing the Overhead



Single instruction multiple Data (SIMD): tens of operations per instruction

I-Cache	RF	Control					
---------	----	---------	--	--	--	--	--

Further specialization (what we achieve using accelerators)							
I-Cache	RF	Control			hundreds	•••	
					or more		

Diagram adapted from W. Qadder, et al., Convolution Engine: Balancing Efficiency & Flexibility in Specialized Computing, ISCA'2013.

Principles for Improving Energy Efficiency

Do less work!

 Amortize overhead in control and data supply across multiple instructions

Do even less work!

 Use smaller (or simpler) data => cheaper operations, lower storage & communication costs

Move data locally and directly

- Store data nearby in simpler memory (e.g., scratchpads are cheaper than cache)
- Wire compute units for direct (or even combinational) communication when possible

Tensor Processing Unit (TPU)

- A domain-specific accelerator specialized for deep learning
 - Main focus: accelerating matrix multiplication (MatMul) with a systolic array
 - Use CISC instructions: MatMul Unit may take thousands of cycles
 - TPUv1 does 8-bit integer (INT8) inference; TPUv2 supports a customized floatingpoint type (bfloat16) for training



Common Hardware Specialization Techniques: A Taxonomy

- Custom Compute Units: Use complex instructions to amortize overhead (e.g., SIMD, "ASIC"-in-an-instruction)
- Custom Numeric Types: Trade off accuracy and efficiency with data types that use smaller bit widths or simpler arithmetic
- Custom Memory Hierarchy: Exploit data access patterns to reduce energy per memory operation
- Custom Communication Architecture: Tailor on-chip networks to data movement patterns

Common Hardware Specialization Techniques: A Taxonomy

- Custom Compute Units: Use complex instructions to amortize overhead (e.g., SIMD, "ASIC"-in-an-instruction)
- Custom Numeric Types: Trade off accuracy and efficiency with data types that use smaller bit widths or simpler arithmetic
- Custom Memory Hierarchy: Exploit data access patterns to reduce energy per memory operation
- Custom Communication Architecture: Tailor on-chip networks to data movement patterns

Customizing Compute Units: An Intuitive View



A simple single-cycle CPU

Evaluating a Simple Expression on CPU



Source: Adapted from Desh Singh's talk at HCP'14 workshop

"Unrolling" the Instruction Execution



Source: Adapted from Desh Singh's talk at HCP'14 workshop

Removing Unused Logic



Source: Adapted from Desh Singh's talk at HCP'14 workshop

An Application-Specific Compute Unit

R5 <= R1 * R3 R6 <= R2 * R4

R7 <= R5 - R6

R8 <= R9 + R7



3. Wire up registers and functional units

Use combinational connections when timing constraints allow (e.g., R7)

Common Hardware Specialization Techniques: A Taxonomy

- Custom Compute Units: Use complex instructions to amortize overhead (e.g., SIMD, "ASIC"-in-an-instruction)
- Custom Numeric Types: Trade off accuracy and efficiency with data types that use smaller bit widths or simpler arithmetic
- Custom Memory Hierarchy: Exploit data access patterns to reduce energy per memory operation
- Custom Communication Architecture: Tailor on-chip networks to data movement patterns

Customized Data Types

Using custom numeric types tailored for a given application/domain improves performance & efficiency



Binary Representation – Positional Encoding

Unsigned number

 MSB has a place value (weight) of 2ⁿ⁻¹

Two's complement

MSB weight = -2ⁿ⁻¹



Fixed-Point Representation of Fractional Numbers

- The positional binary encoding can also represent fractional values, by using a **fixed** position of the binary point and place values with negative exponents
 - (-) Less convenient to use in software, compared to floating point
 - (+) Much more efficient in hardware



Next Lecture

More Hardware Specialization

Acknowledgements

- These slides contain/adapt materials developed by
 - Bill Dally, NVIDIA
 - System for AI Education Resource by Microsoft Research