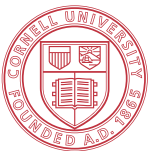




ECE 6775  
High-Level Digital Design Automation  
Fall 2025

# Hardware Specialization



Cornell University



# Announcements

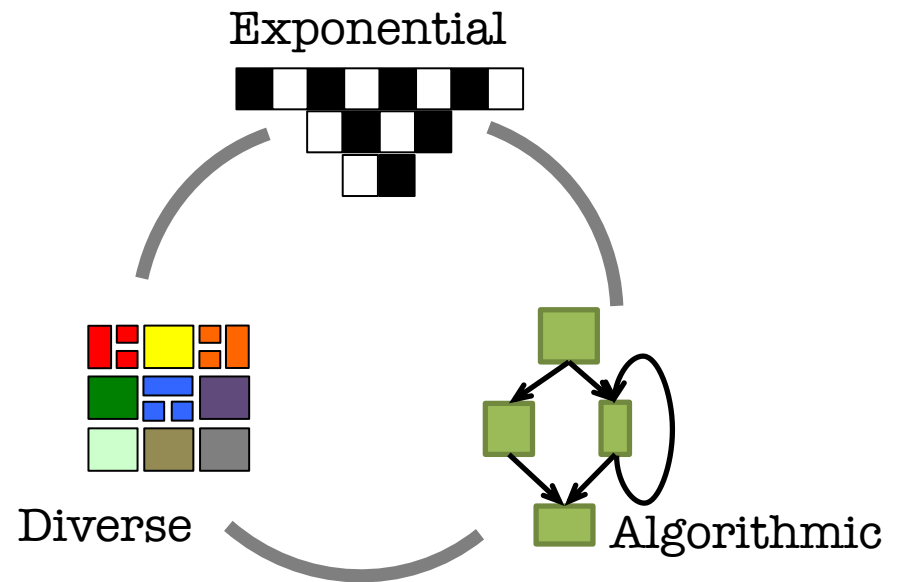
- ▶ Vivado HLS setup guide will be released soon
  - Complete tool setup before Thursday 9/4
- ▶ First paper reading session on Tuesday 9/9
  - A. Boutros and V. Betz, “[FPGA Architecture: Principles and Progression](#)”, IEEE CAS-M 2021
  - Complete reading before class

# Recap: Our Interpretation of E-D-A

**Exponential**  
in complexity (or **Extreme** scale)

**Diverse**  
increasing system heterogeneity

**Algorithmic**  
intrinsically computational

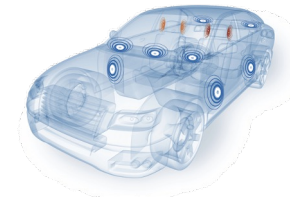
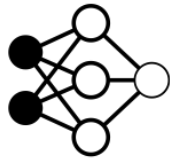


# Agenda

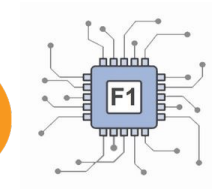
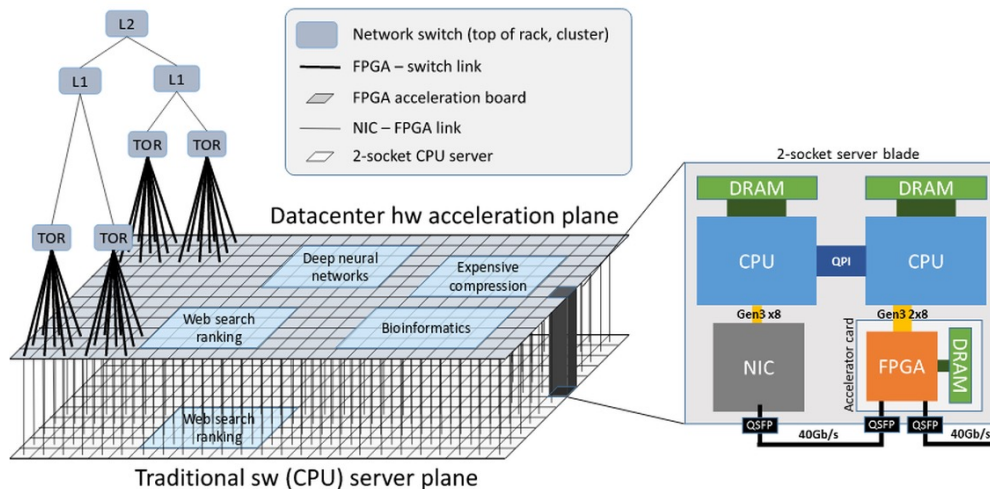
- ▶ Motivation for hardware specialization
  - Key driving forces from applications and technology
  - Main sources of inefficiency in general-purpose computing
- ▶ A taxonomy of common specialization techniques
- ▶ Introduction to fixed-point types

# A Golden Age of Hardware Specialization

- ▶ **Higher demand** on efficient compute acceleration, esp. for AI workloads



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



# Rising Computational Demands of Emerging Applications

- ▶ Deep neural networks (DNNs) require enormous amount of compute
  - Consider ResNet50, a 50-layer model that performs ~4 billion operations to classify one image (a small model by today's standards)

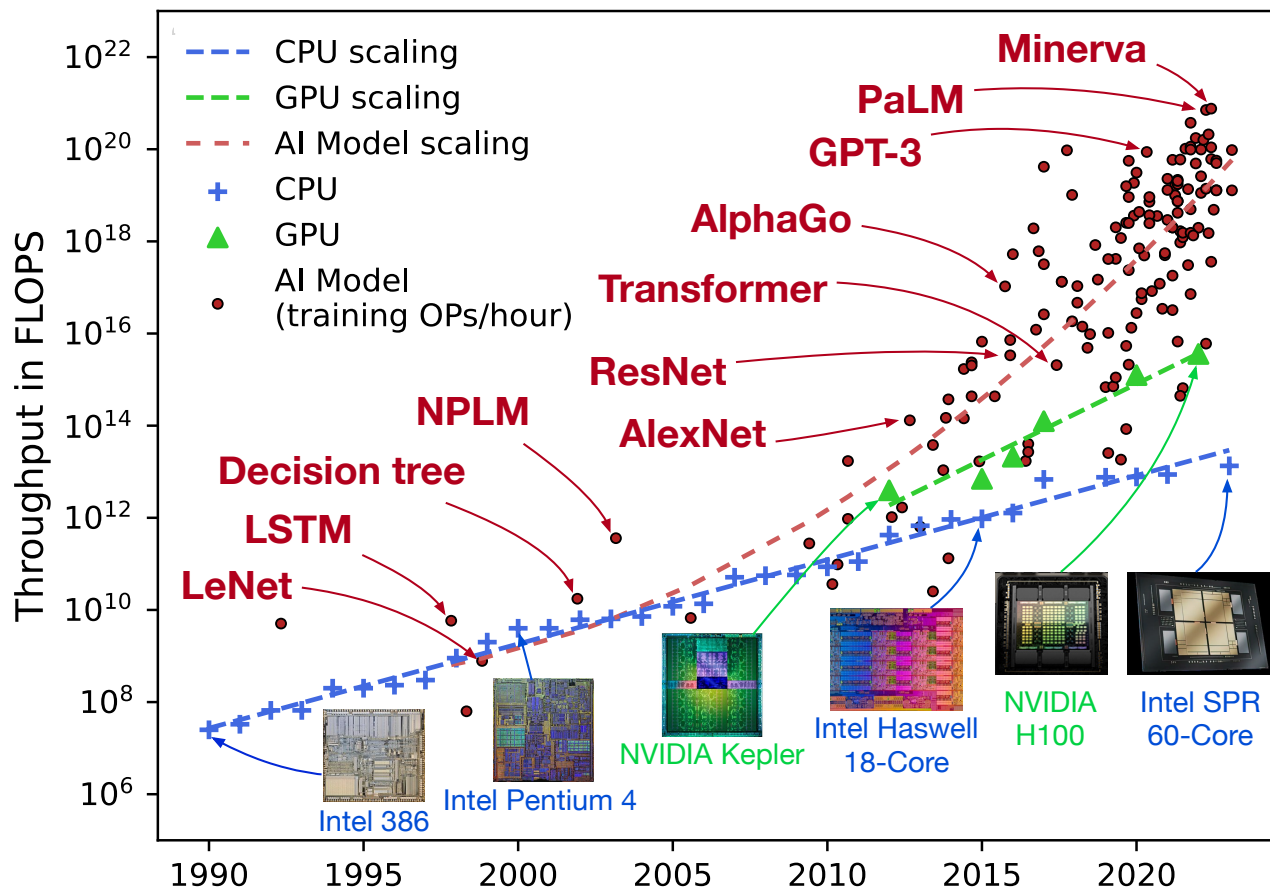
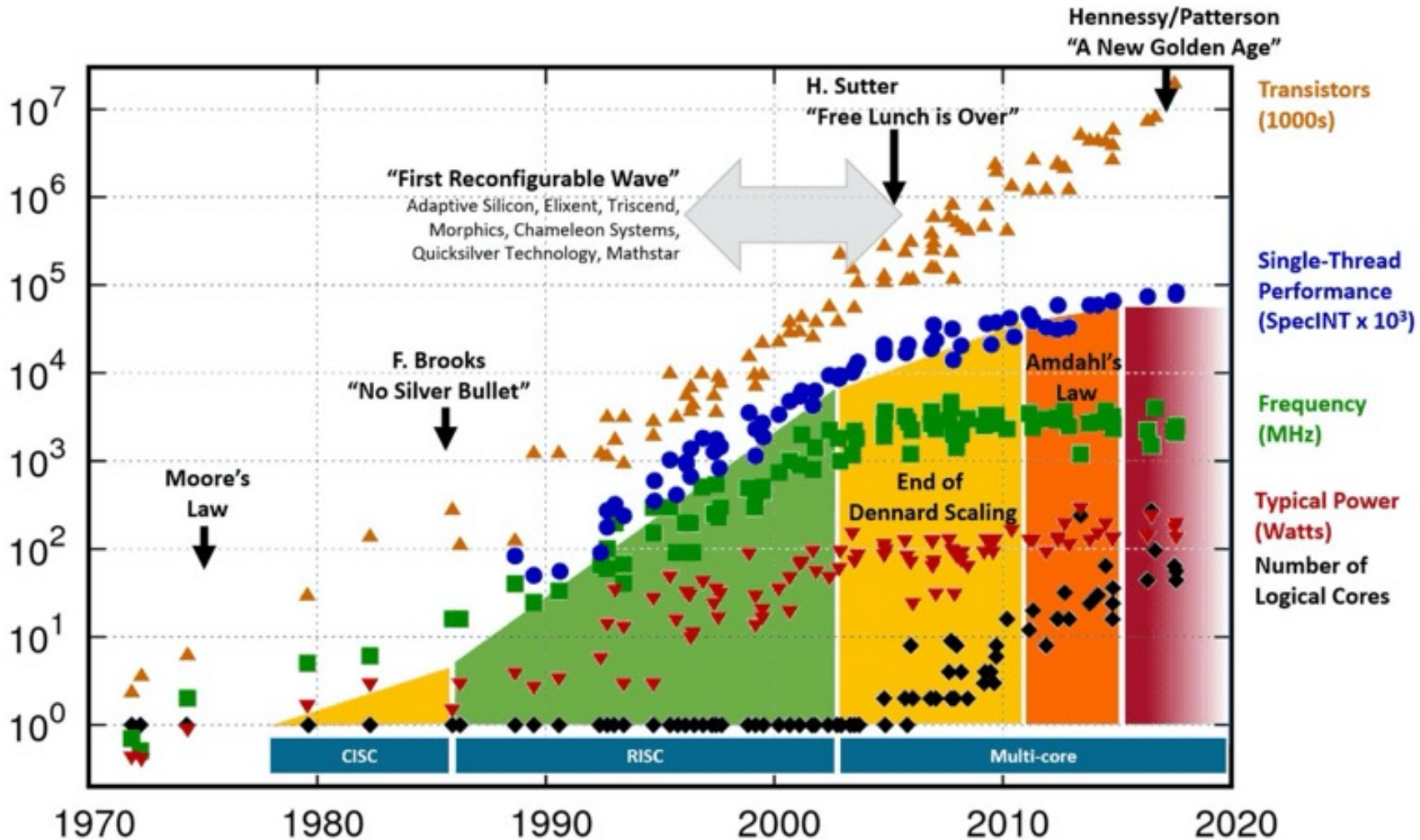


Figure source: Cornell Zhang Research Group

# 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.karlsruhp.net/2018/02/42-years-of-microprocessor-trend-data/>; “First Wave” added by Les Wilson, Frank Schirrmester  
 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** (S: the scaling factor)

Transistor (T) count	$S^2$
Capacitance / T	$1/S$
Voltage ( $V_{dd}$ )	$1/S$
Frequency	$S$
<b>Total power</b>	<b>1</b>

**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

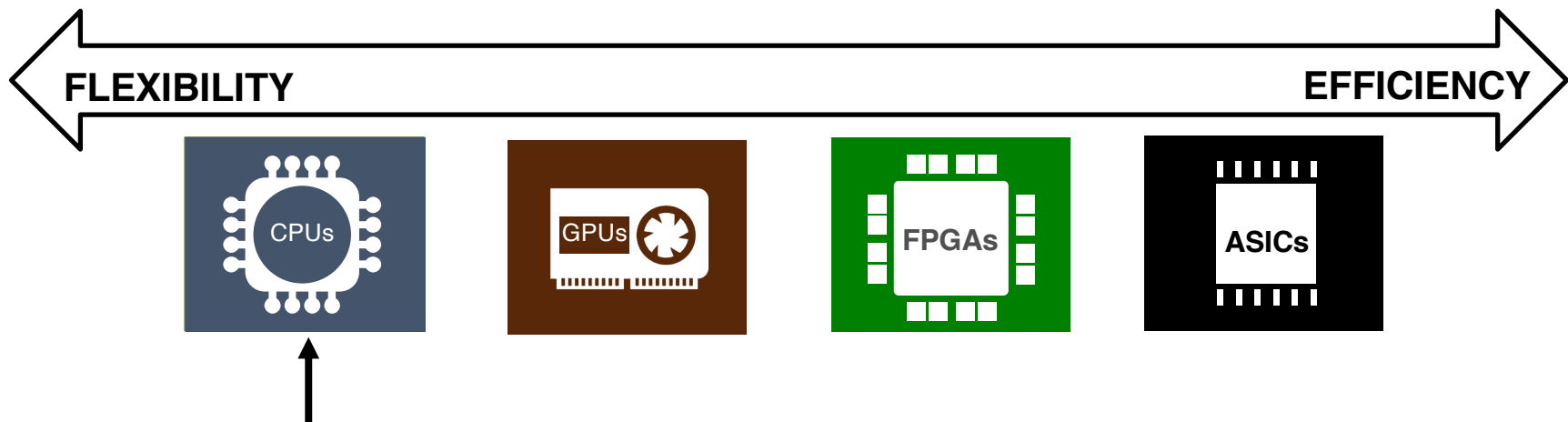
## Leakage limited scaling

Transistor (T) count	$S^2$
Capacitance / T	$1/S$
Voltage ( $V_{dd}$ )	$\sim 1$
Frequency	$\sim 1$
<b>Total power</b>	<b>S</b>

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

- ▶ Implication: “Dark silicon”?
  - Power limits restrict how much of the chip (not 100%) can be activated simultaneously
- ▶ **Recent trend:** Datacenter GPUs are breaking past the “flat” power envelope via liquid cooling & advanced power delivery
  - Hopper H100:  $\sim 700W$   $\rightarrow$  Blackwell B200:  $\sim 1200W$

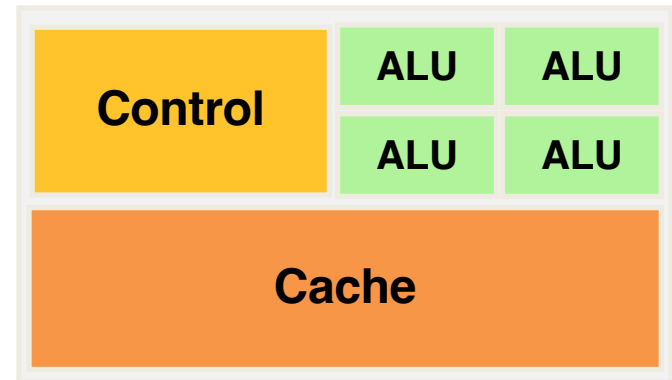
# Trade-off Between Flexibility and Efficiency



Why are general-purpose CPUs less energy efficient?

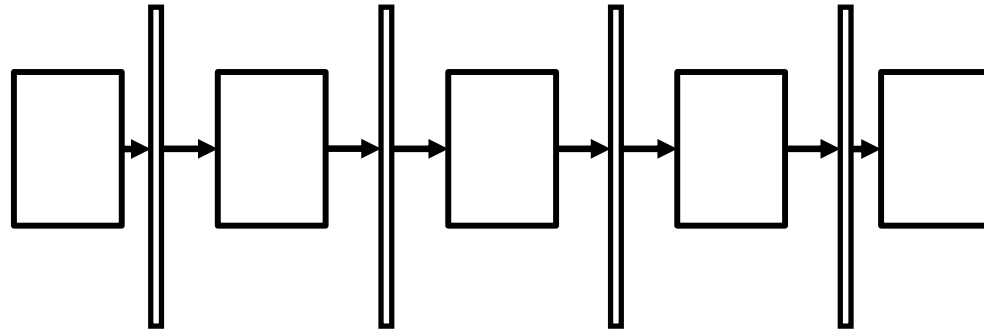
# CPU Core Architecture

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

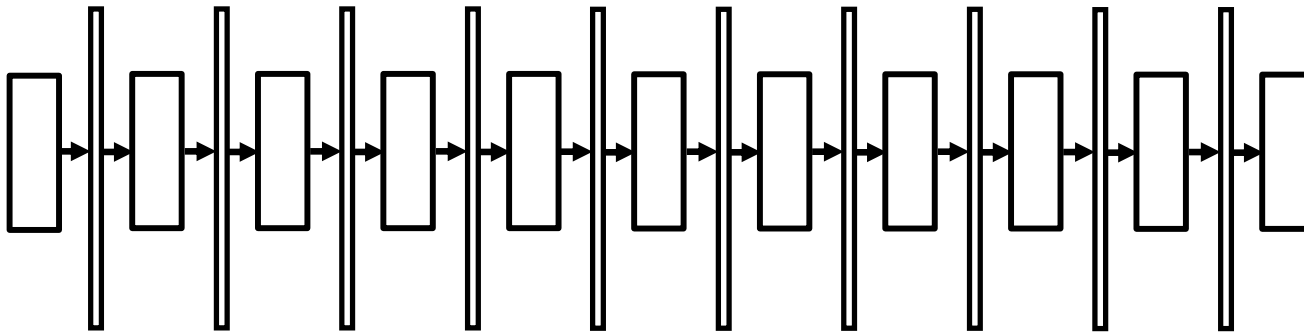


- ▶ **Mainly optimized to reduce latency of running serial code**
  - Shallow pipelines (< 30 stages)
  - Superscalar, out-of-order (OOO) execution, speculative execution, branch prediction, prefetching, etc.
  - **Low throughput, even with multithreading**

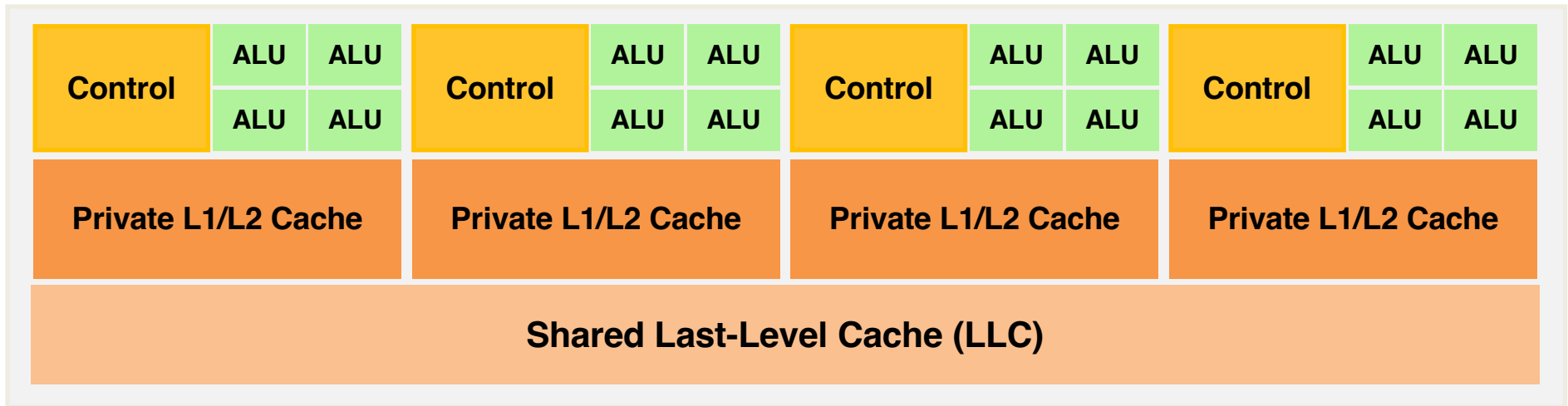
# Discussion



Shallow vs. Deep Pipelining in context of CPU design



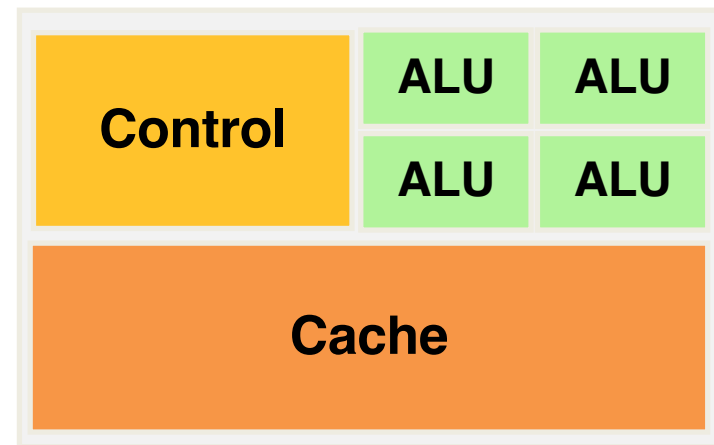
# Multi-Core Architecture



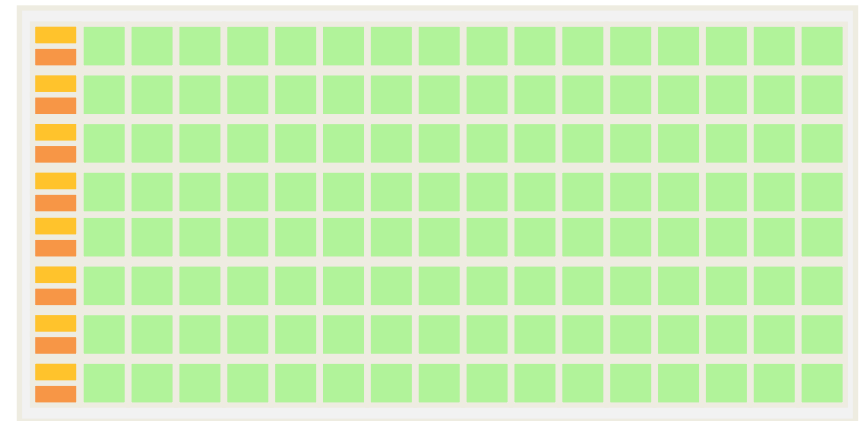
With four cores, should we expect a 4x speedup on an arbitrary application?

# 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 data-parallel workloads**
  - Initially targeting graphics code
  - Latency tolerant with many concurrent threads



**CPU**



**GPU**

# 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
Small SRAM (8KB)	5 pJ
Moving two 32b words by	Energy
40mm (across a 400mm <sup>2</sup> chip)	77 pJ
1mm (local communication)	1.9 pJ
Arithmetic on two 32b words	Energy
FMA (float fused multiply-add)	1.2 pJ
IADD (integer add)	0.02 pJ

**65,000x** (Red arrow from DRAM to IADD)

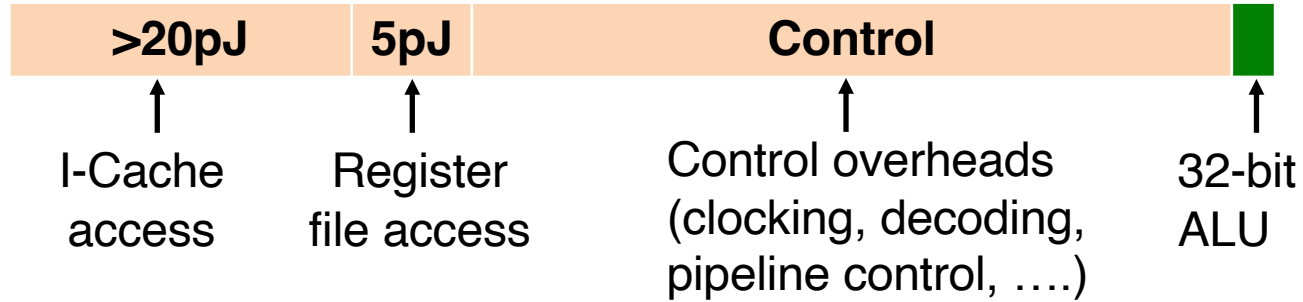
**250x** (Orange arrow from FMA to IADD)

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



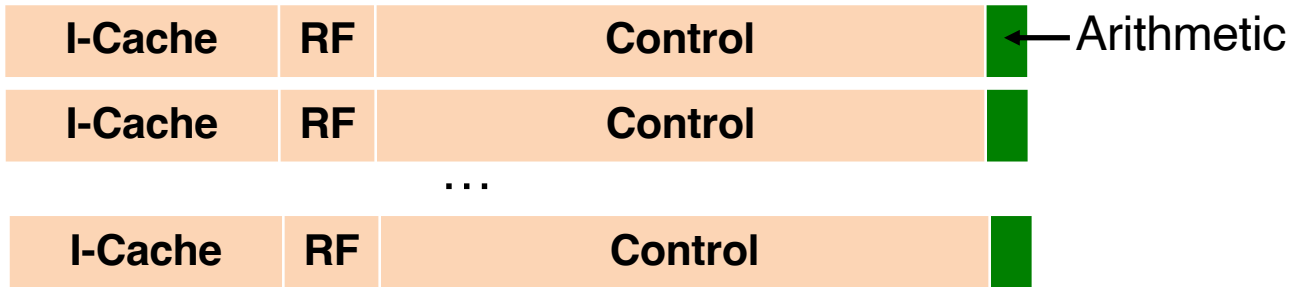
# Principles for Improving Energy Efficiency

Do less work!

- **Amortize overhead** in control and data supply across multiple instructions

# Amortizing the Overhead

A sequence of energy-inefficient instructions



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



Further specialization (what we achieve using accelerators)



# Principles for Improving Energy Efficiency

Do less work!

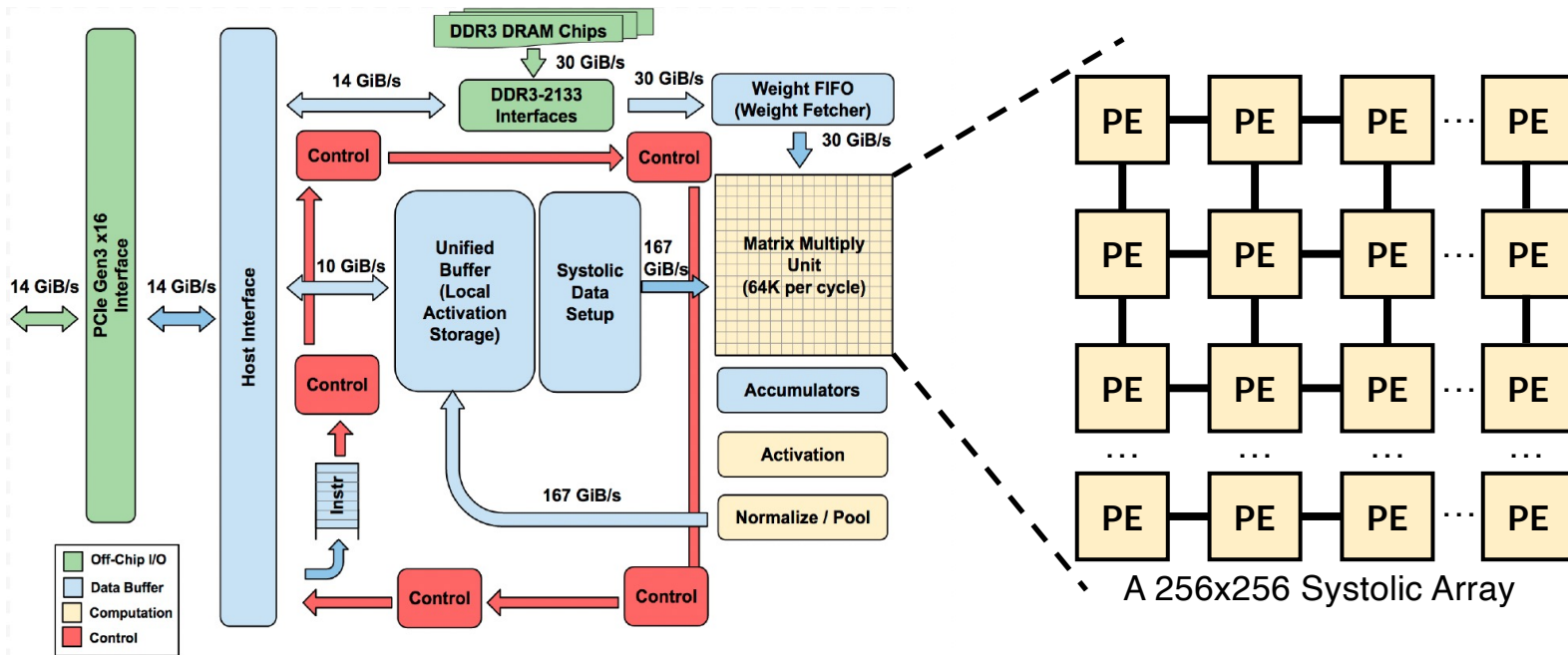
- **Amortize overhead** in control and data supply across multiple instructions

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

# Example: Tensor Processing Unit (TPU)

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



Google TPU v1

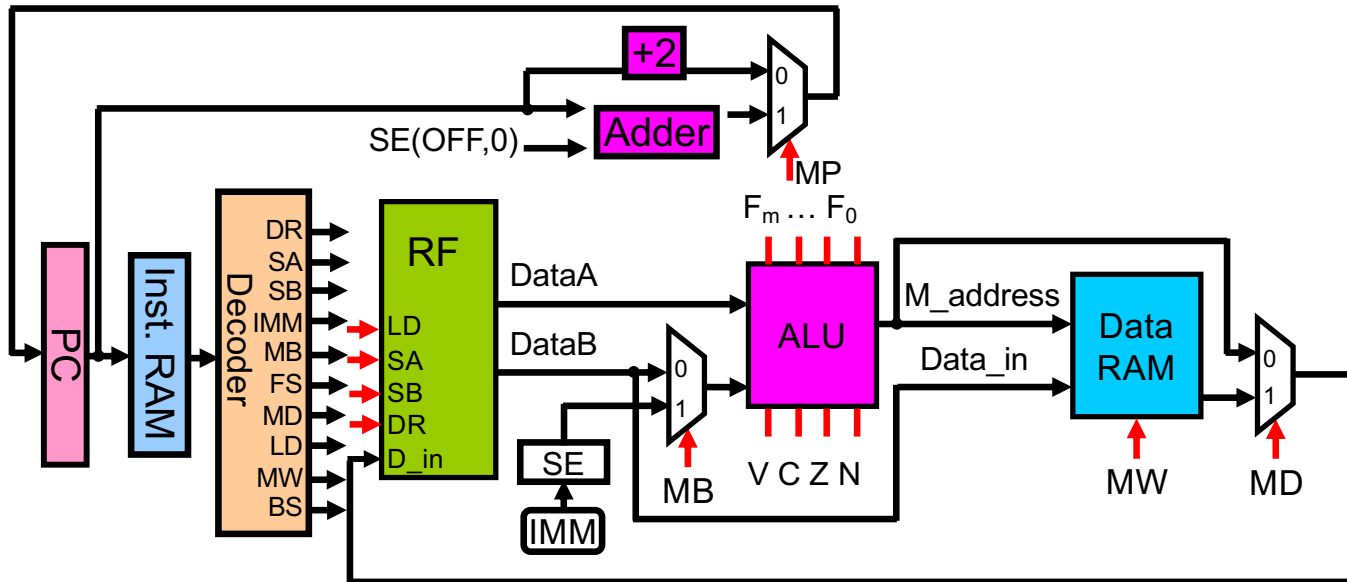
# Common Hardware Specialization Techniques

- ▶ **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

- ▶ **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

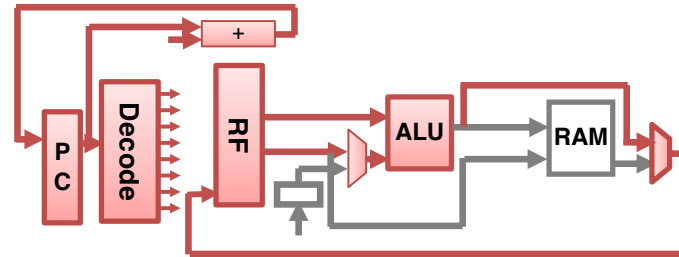
# Customizing Compute Units: An Intuitive View



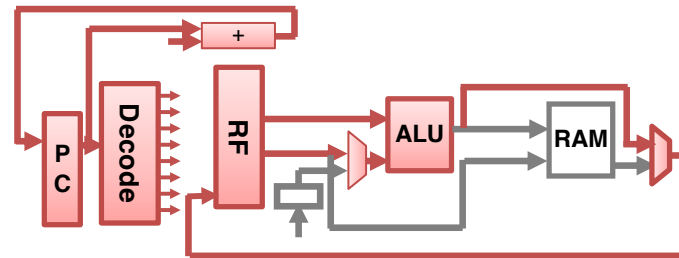
A simple single-cycle CPU

# Evaluating a Simple Expression on CPU

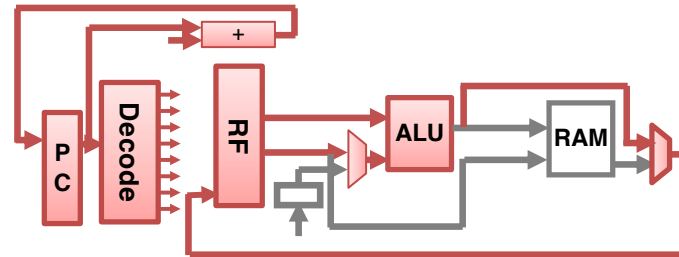
$$R5 \leftarrow R1 * R3$$



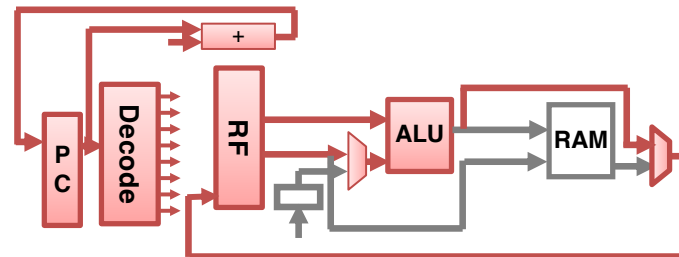
$$R6 \leftarrow R2 * R4$$



$$R7 \leftarrow R5 - R6$$



$$R8 \leftarrow R9 + R7$$

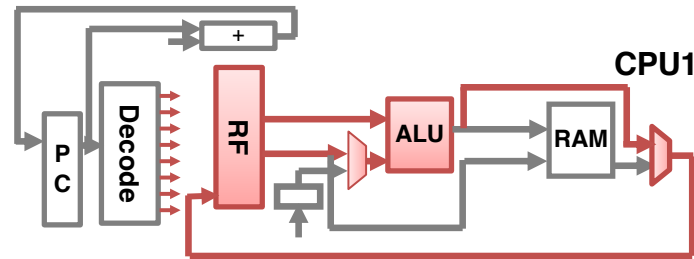


Cycle-by-cycle  
CPU activities

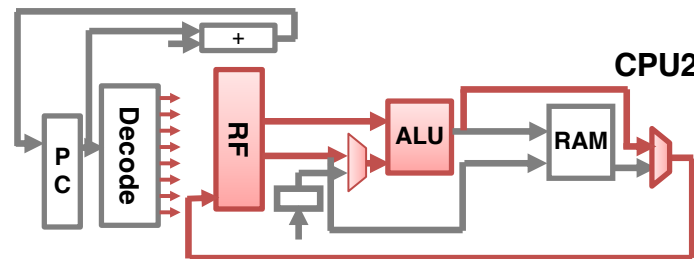


# “Unrolling” the Instruction Execution

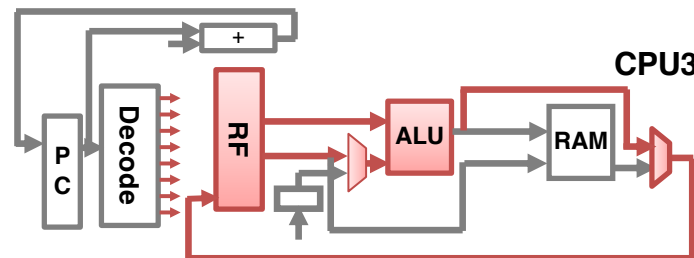
$R5 \leftarrow R1 * R3$



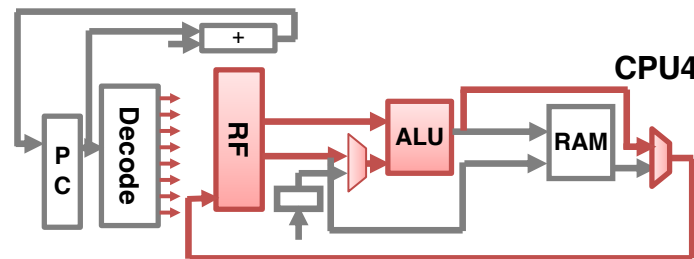
$R6 \leftarrow R2 * R4$



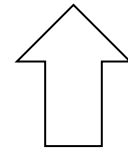
$R7 \leftarrow R5 - R6$



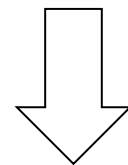
$R8 \leftarrow R9 + R7$



1. Replicate the CPU hardware  
Instruction fixed  
=> disable fetch  
& decode logic

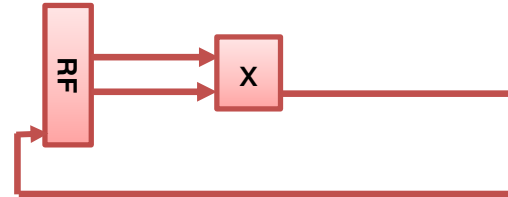


Space

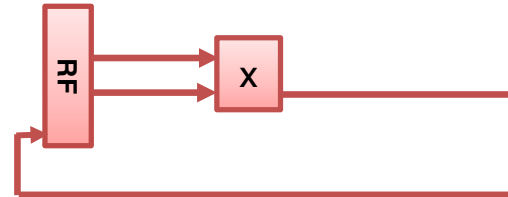


# Removing Unused Logic

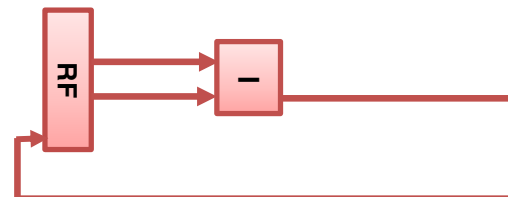
$$R5 \leftarrow R1 * R3$$



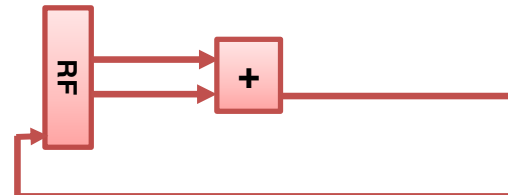
$$R6 \leftarrow R2 * R4$$



$$R7 \leftarrow R5 - R6$$



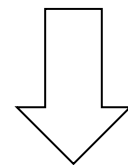
$$R8 \leftarrow R9 + R7$$



**2. Removing unused logic**  
=> ALU also simplified



**Space**



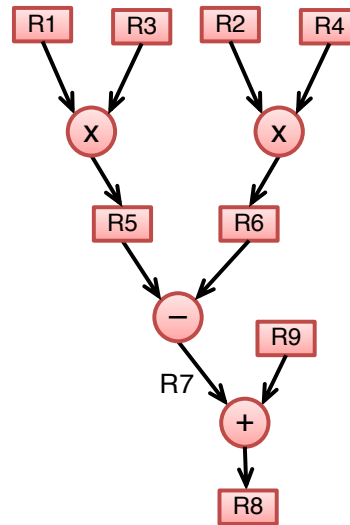
# An Application-Specific Compute Unit

$$R5 \leftarrow R1 * R3$$

$$R6 \leftarrow R2 * R4$$

$$R7 \leftarrow R5 - R6$$

$$R8 \leftarrow R9 + R7$$



### 3. Wire up registers and functional units

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



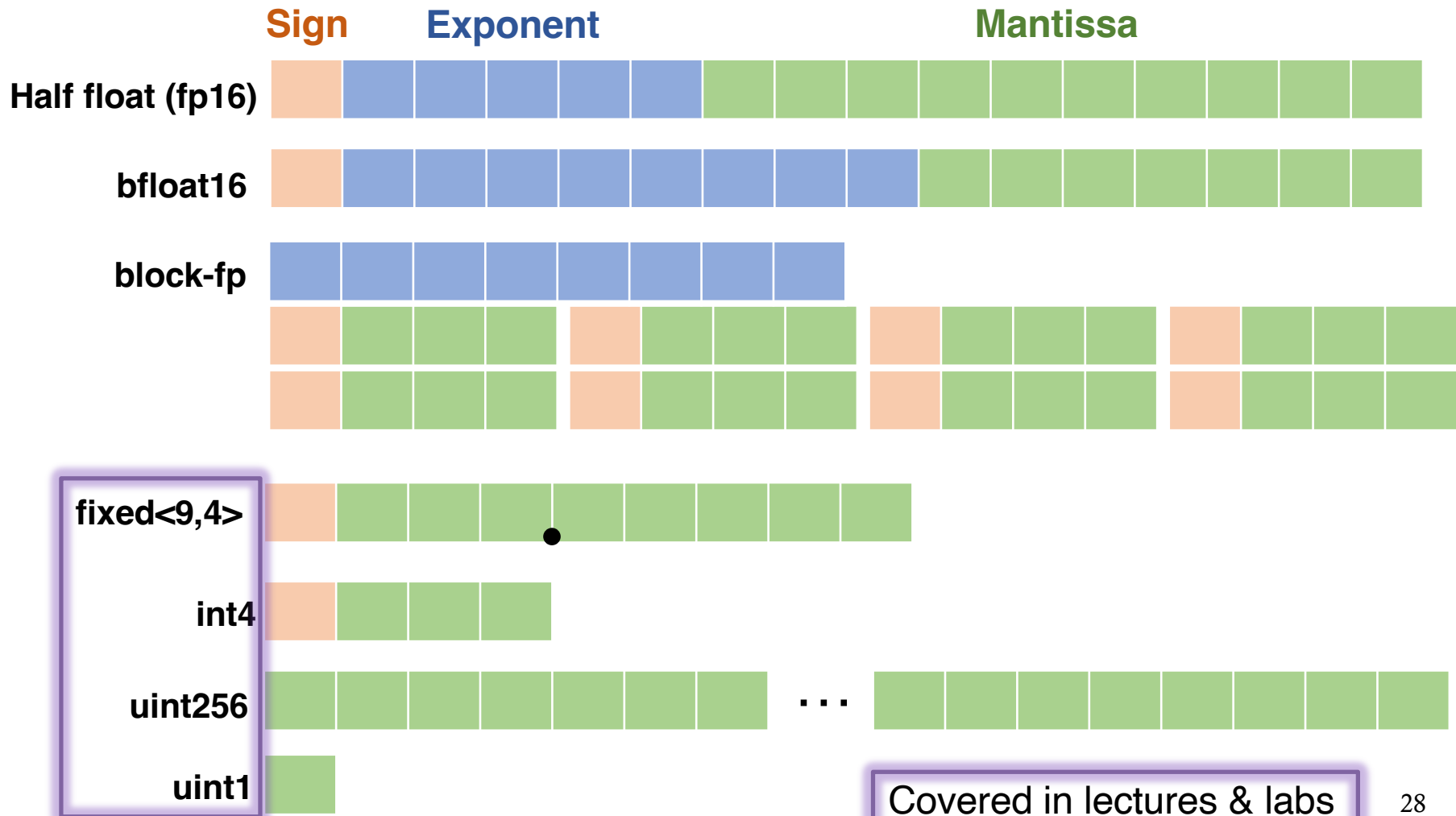
“ASIC”-in-an-instruction

## 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



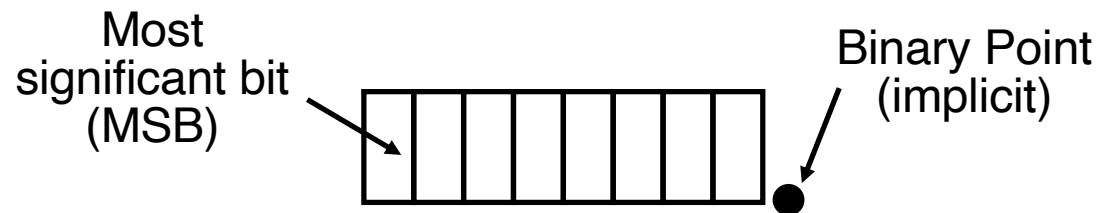
# Binary Representation – Positional Encoding

## Unsigned number

- ▶ MSB has a place value (weight) of  $2^{n-1}$

## Two's complement

- ▶ MSB weight =  $-2^{n-1}$

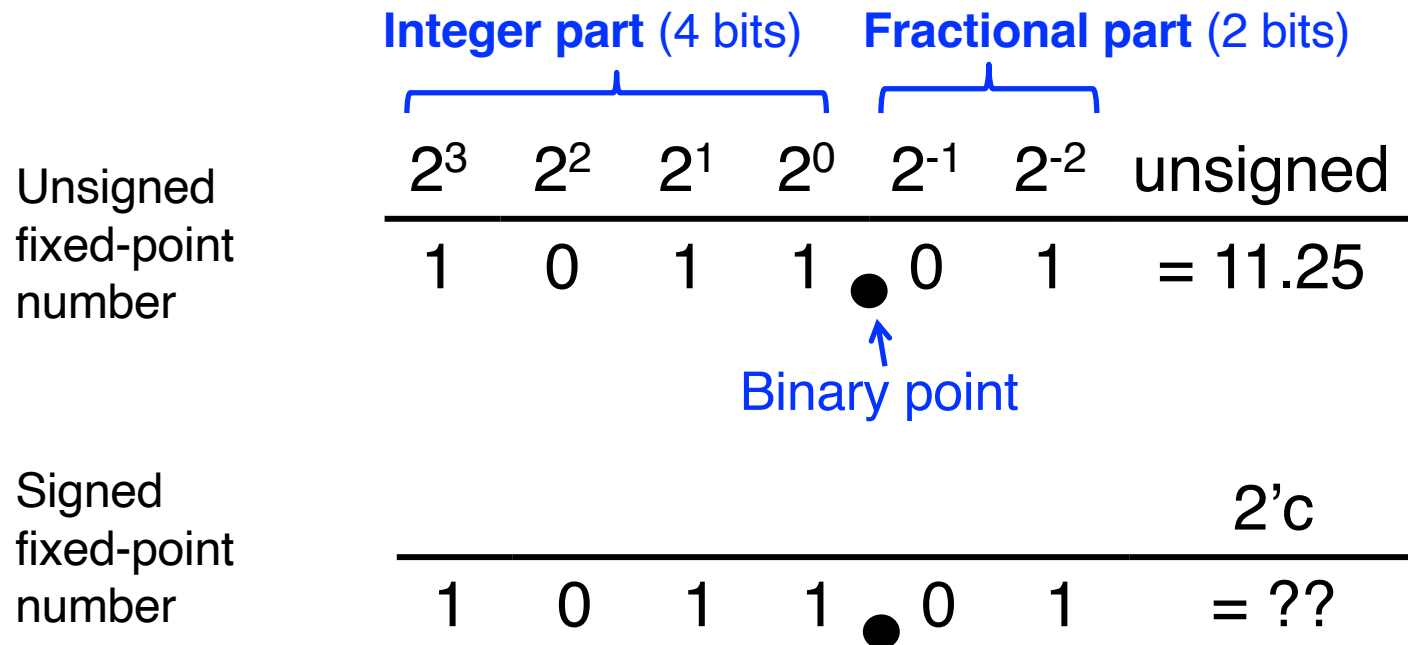


$2^3$	$2^2$	$2^1$	$2^0$	unsigned
1	0	1	1	= 11

$-2^3$	$2^2$	$2^1$	$2^0$	2's c
1	0	1	1	= -5

# 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