

Principles of Computer Architecture

CSE 240A

Fall 2024

Hadi Esmaeilzadeh

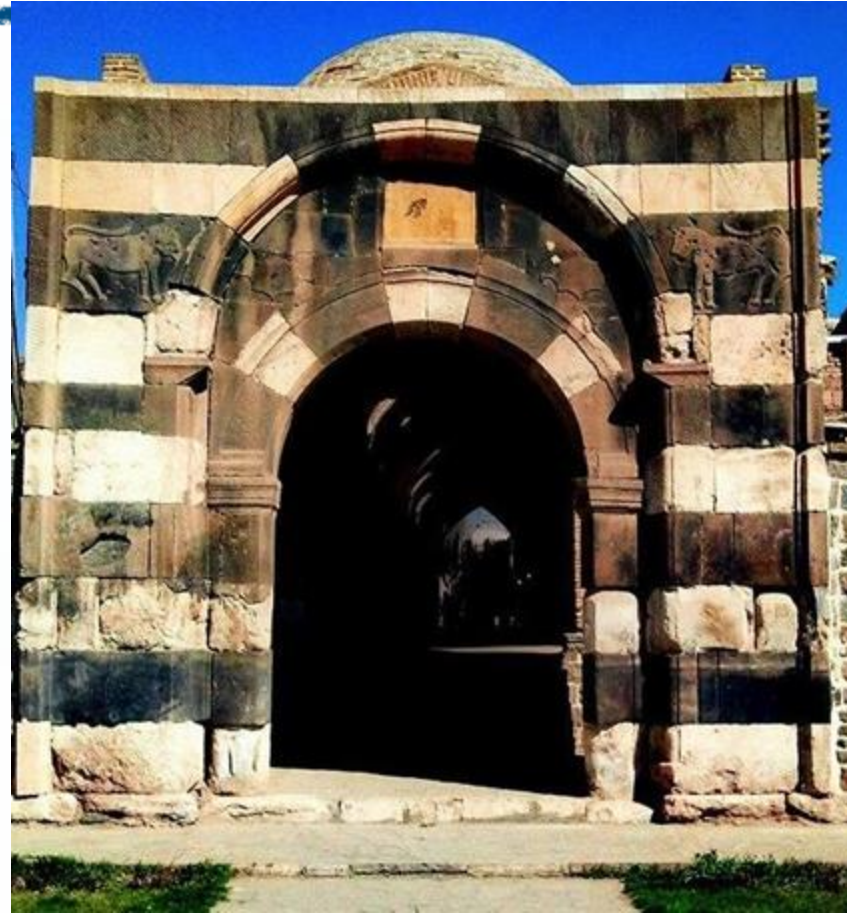
hadi@ucsd.edu

University of California, San Diego



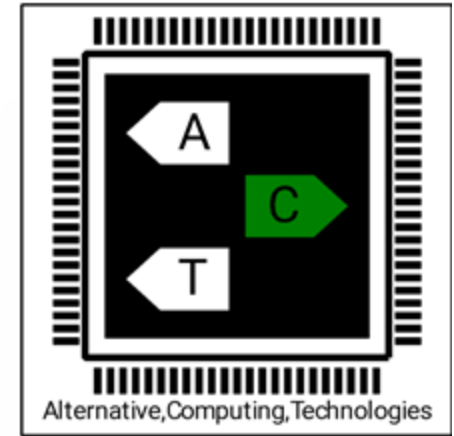
Hadi Esmaeilzadeh

From Khoy, Iran



Research: ACT Lab

Alternative Computing Technologies



- System design for immersive machine intelligence
- System design for robotics
- System design for health
- Analog computing
- General-purpose approximate computing
- Bridging neuromorphic and von Neumann models of computing

Agenda

1. Who is Hadi

2. Course organization

3. Why CS 240A Principles of Computer Architecture

Objective



- Learn more about the fundamental design tradeoffs in computer architecture and some of the recent research issues/trends.
- To provide the necessary background and experience to do design and research in computer system design.
- Strong emphasis on
 - Quantitative Evaluations and Trade-offs
 - Hands on programming assignments

Course Information

- Course Resources
 - Gradescope: All assignment submissions
 - Piazza: Course Resources and Q/A
 - TA: Hanyang Xu, hanyang@ucsd.edu
 - Office Hour: 930am-1030am CSE 3rd floor lobby

Format



- **Lectures** are the main source for exams and homework
- There is no perfect textbook for this course!
 - Recommended reading:
Sixth Edition of Computer Architecture: A Quantitative Approach by John Hennessy and David Patterson **AND** Microprocessor architecture by Jean-Loup Baer, Cambridge
- Read the related papers, and do the programming assignments

Grading rubric



Component	Fraction
Class Participation	5%
Scribes	10%
Homeworks	15%
Midterm	25%
Projects	25%
Final Exam	40%
Total	120%

Project Assignments

- This course requires heavy programming
- Don't take too many program/project heavy courses together!
- It is 4-credit course but you feel a 5-6 credit course
- The most CSlike course in ECE, the most ECElike course in CS
- Each individual should be expert in all aspects of the work
- Please **DO NOT CHEAT!** It is just not cool!
 - Follow the UCSD Academic Honor Code
 - Ask me if you are not sure



Agenda

1. Who is Hadi

2. Course organization

3. Why CSE 240A

1. How we became an industry of new capabilities

2. Why we might become an industry of replacement

What has made computing pervasive? What is the backbone of computing industry?



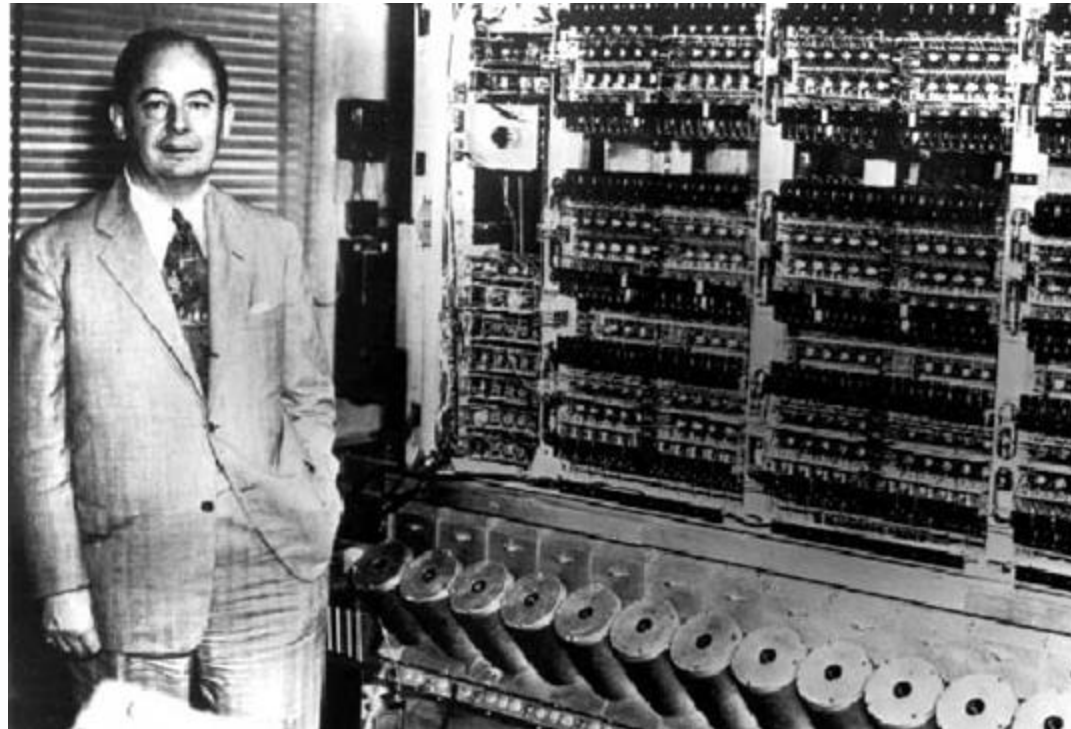
Programmability

```
public class TcpClientSample
{
    public static void Main()
    {
        byte[] data = new byte[1024]; string input, stringData;
        TcpClient server;
        try{
            server = new TcpClient(" . . . . ", port);
        }catch (SocketException){
            Console.WriteLine("Unable to connect to server");
            return;
        }
        NetworkStream ns = server.GetStream();
        int recv = ns.Read(data, 0, data.Length);
        stringData = Encoding.
            ASCII.GetString(data, 0, recv);
        Console.WriteLine(stringData);
        while(true){
            input = Console.ReadLine();
            if (input == "exit") break;
            newchild.Properties["ou"].Add
                ("Auditing Department");
            newchild.CommitChanges();
            newchild.Close();
        }
    }
}
```

Networking



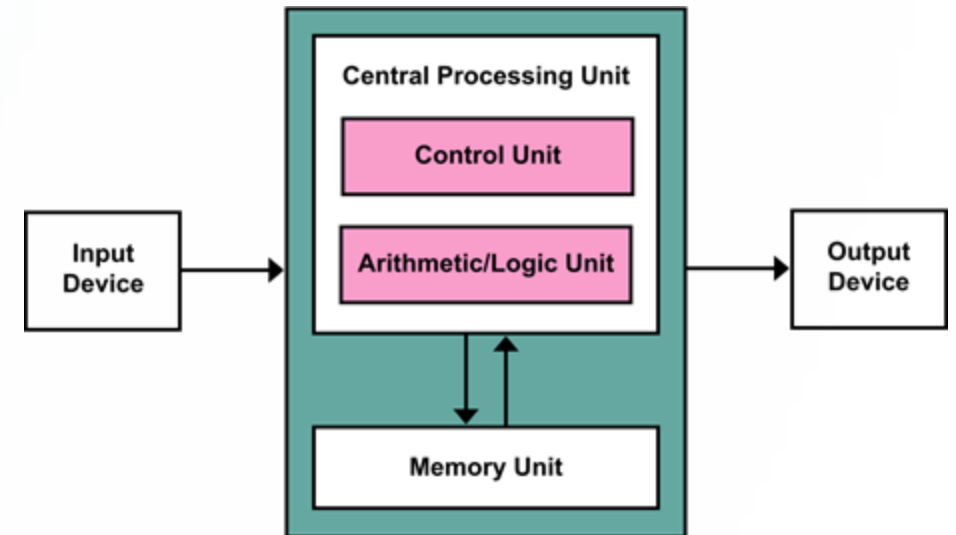
What makes computers programmable?



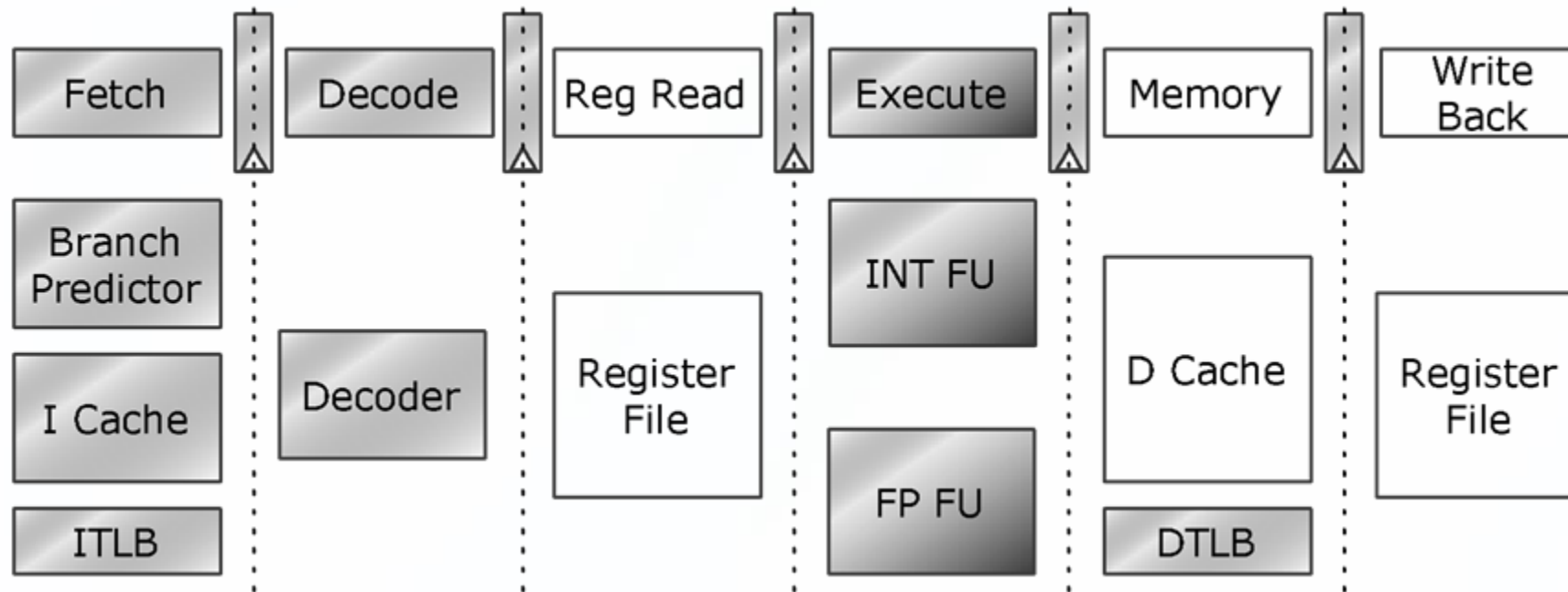
von Neumann architecture

General-purpose processors

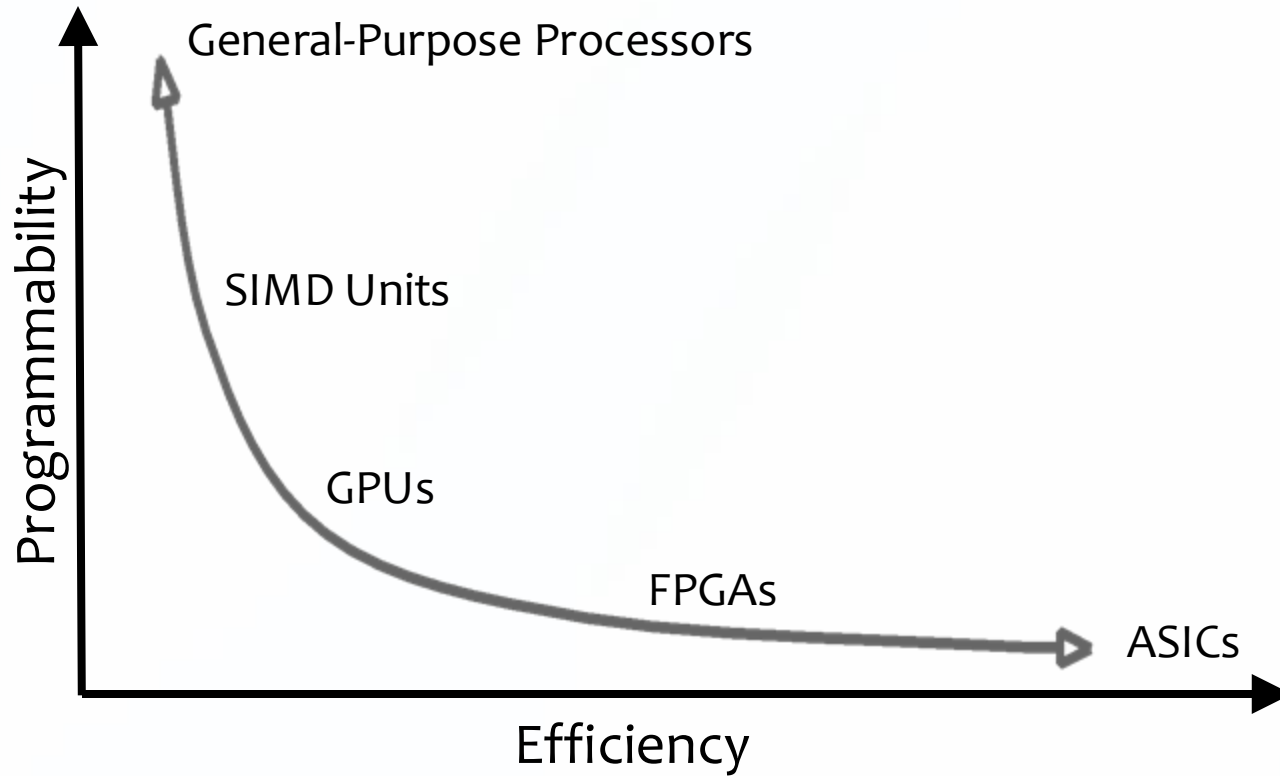
- Components
 - Memory (RAM)
 - Central processing unit (CPU)
 - Control unit
 - Arithmetic logic unit (ALU)
 - Input/output system
- Memory stores program and data
- Program instructions execute sequentially



Programmability versus Efficiency



Programmability versus Efficiency



What is the difference between the computing industry and the toothpaste industry?



Industry of Replacement



1971

2024



Industry of New Capabilities

Can we continue being an industry of new possibilities?

Personalized
healthcare

Virtual
reality

Real-time
translators

Agenda

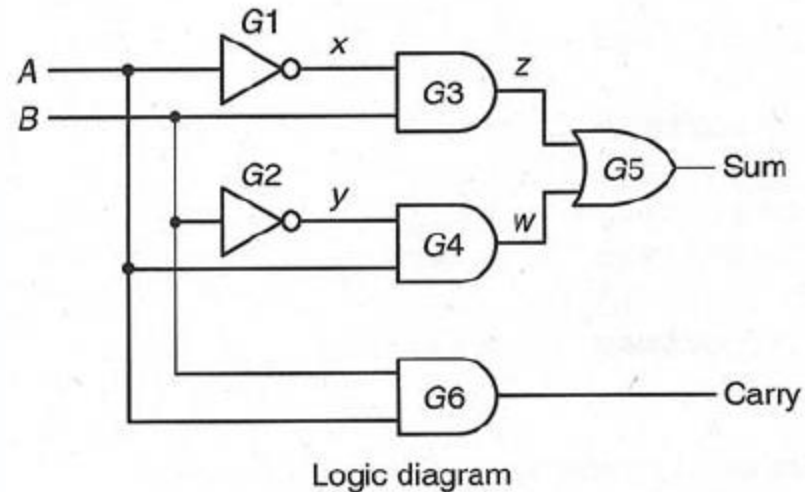
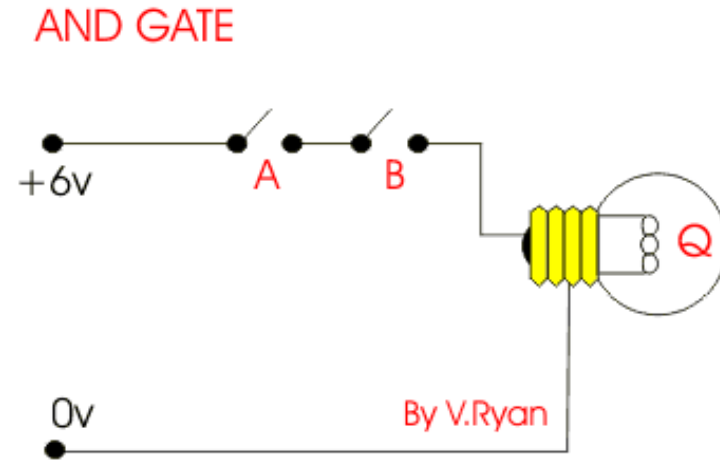
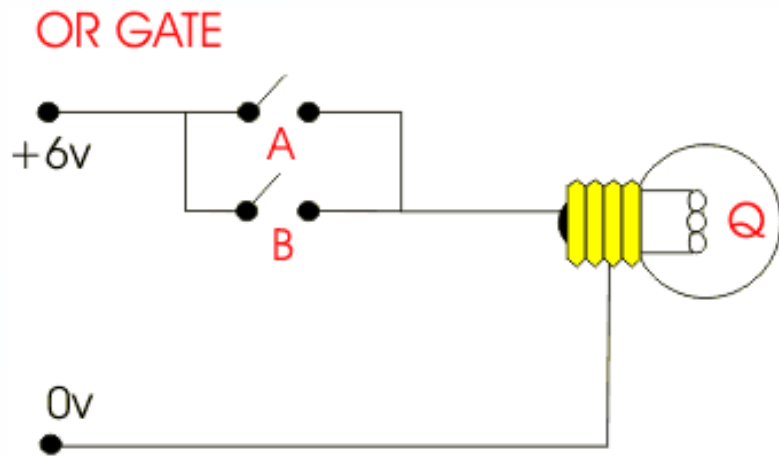
1. Who is Hadi
2. Course organization
3. Why CSE 240C Advanced Microarchitecture

1. How we became and industry of new capabilities

2. Why we might become an industry of replacement
3. Specialization and accelerators

Transistors/switches

Building blocks of computing



Moore's Law

Or, how we became an industry of new capabilities

Every 2 Years

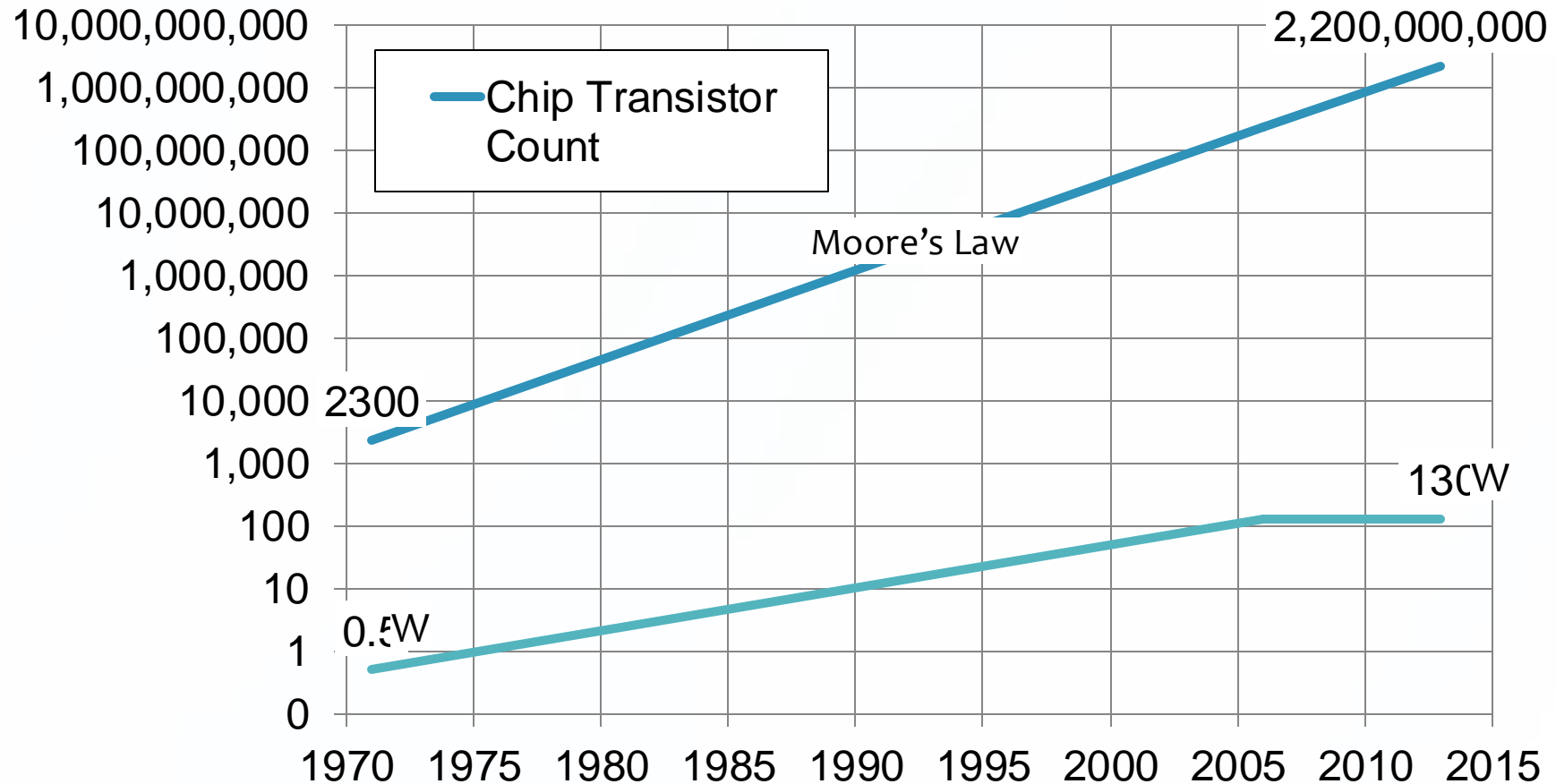
Double the number of transistors

Build higher performance
general-purpose processors

- Make the transistors available to masses
- Increase performance ($1.8\times\uparrow$)
- Lower the cost of computing ($1.8\times\downarrow$)

What is the catch?

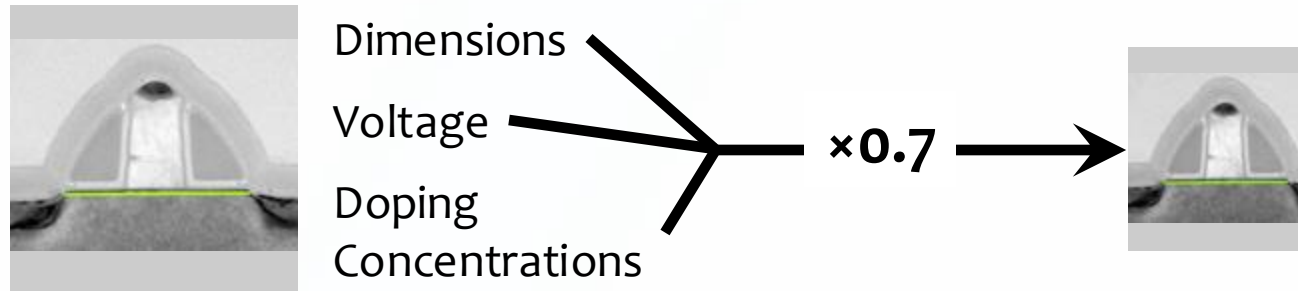
Powering the transistors without melting the chip



Dennard scaling:

Doubling the transistors; scale their power down

Transistor: 2D Voltage-Controlled Switch



Area $\xrightarrow{0.5\times\downarrow}$

Capacitance $\xrightarrow{0.7\times\downarrow}$

Frequency $\xrightarrow{1.4\times\uparrow}$

$$\text{Power} = \text{Capacitance} \times \text{Frequency} \times \text{Voltage}^2$$

Power $\xrightarrow{0.5\times\downarrow}$

Dennard Scaling Broke:

~~Double the transistors; still scale their power down~~

Transistor: 2D Voltage-Controlled Switch



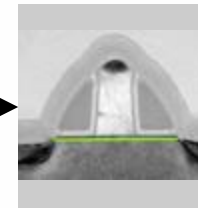
Dimensions

~~Voltage~~

Doping

Concentrations

$\times 0.7$



Area $\xrightarrow{0.5\times\downarrow}$

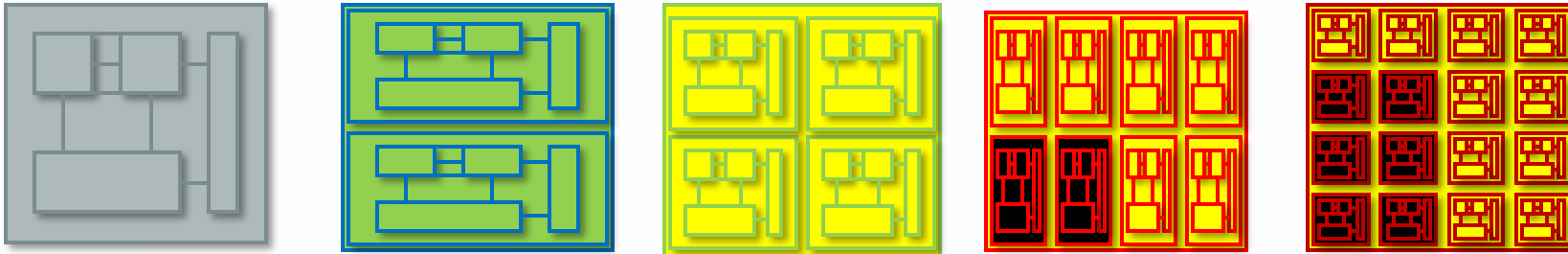
Capacitance $\xrightarrow{0.7\times\downarrow}$

Frequency $\xrightarrow{1.4\times\uparrow}$

$$\text{Power} = \text{Capacitance} \times \text{Frequency} \times \text{Voltage}^2$$

Power $\xrightarrow{0.5\times\downarrow}$

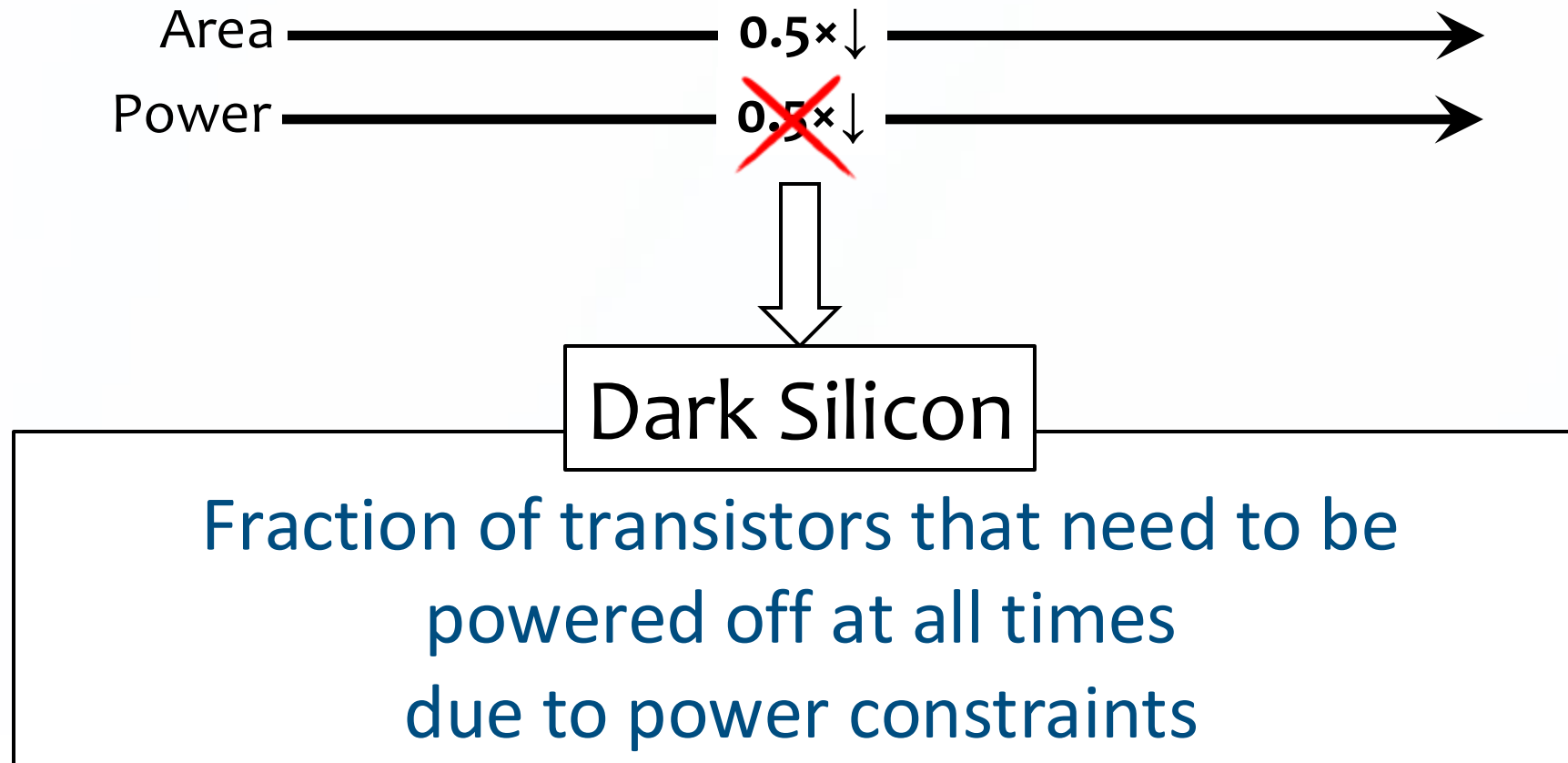
Why Diminishing Returns?



- Transistor area is still scaling
- Voltage and capacitance scaling have slowed
- Result: designs are power, not area, limited

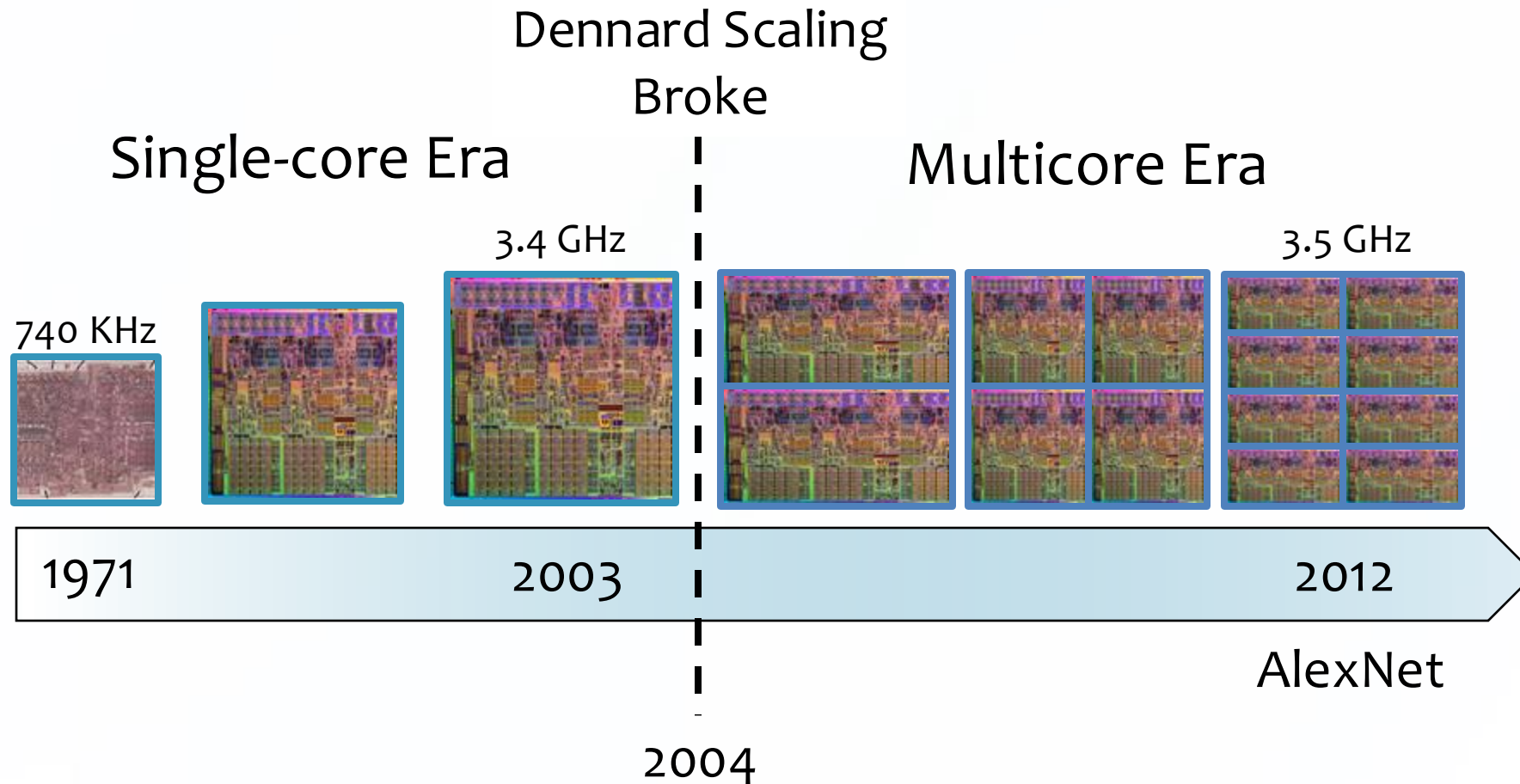
Dark Silicon

If you cannot power them, why bother making them?



Looking Back

Evolution of processors



Are multicores a long-term solution or just a stopgap?

Agenda

1. Who is Hadi
2. Course organization
3. Why CSE 240C Advanced Microarchitecture
 1. How we became and industry of new capabilities
 - 2. Why we might become an industry of replacement**
 3. Specialization and accelerators

Modeling future multicores

Quantify the severity of the problem

Predict the performance of best-case multicores

- From 45 nm to 8 nm
- Parallel benchmarks
- Fixed power and area budget

**Transistor
Scaling Model**

**Single-Core
Scaling Model**

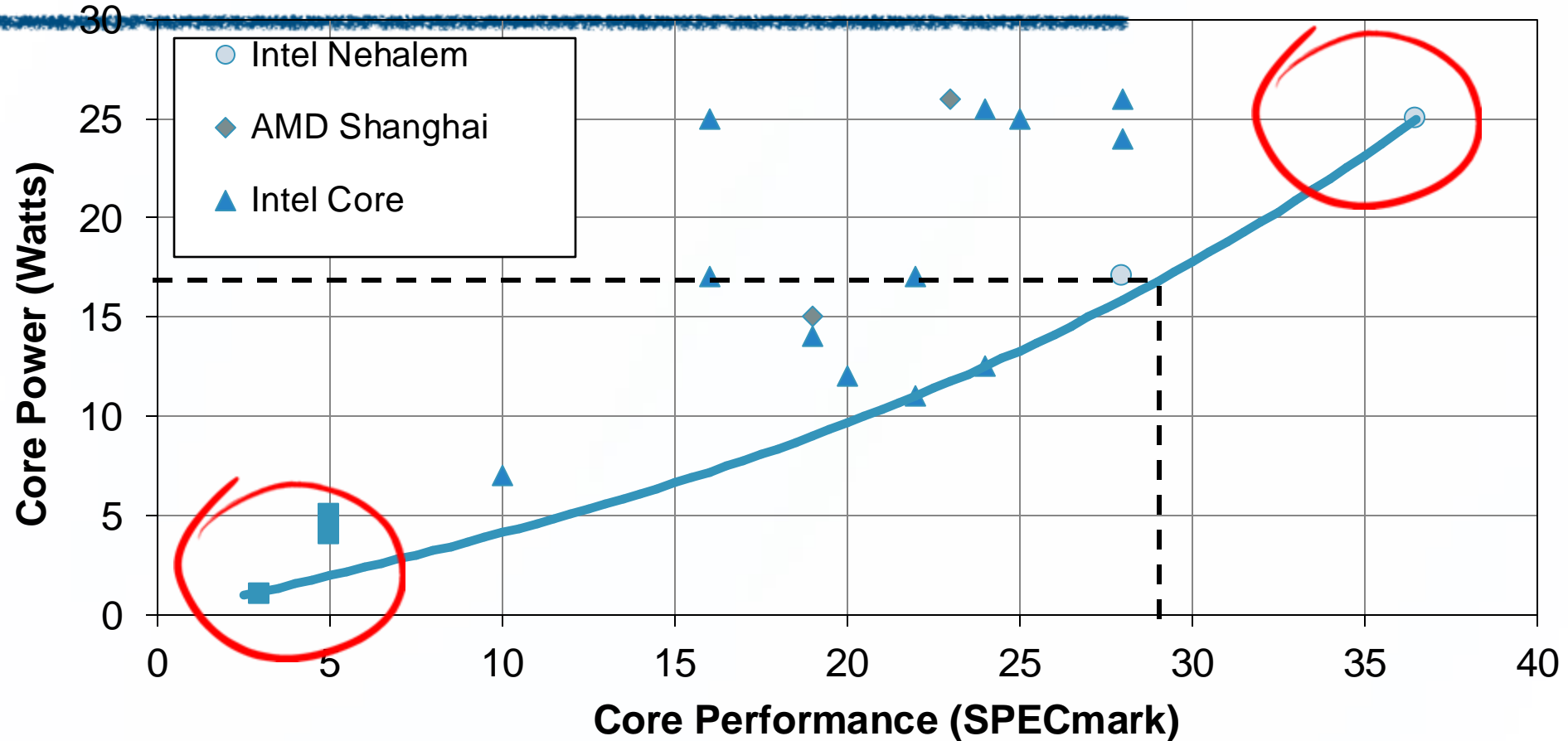
**Multicore
Scaling Model**

Transistor scaling model

From 45 nm to 8 nm

	[Dennard, 1974]	[ITRS, 2010]	[VLSI-DAT, 2010]
	Historical Scaling	Optimistic Scaling Model	Conservative Scaling Model
Area	32× ↓	32× ↓	32× ↓
Power	32× ↓	8.3× ↓	4.5× ↓
Speed	5.7× ↑	3.9× ↑	1.3× ↑

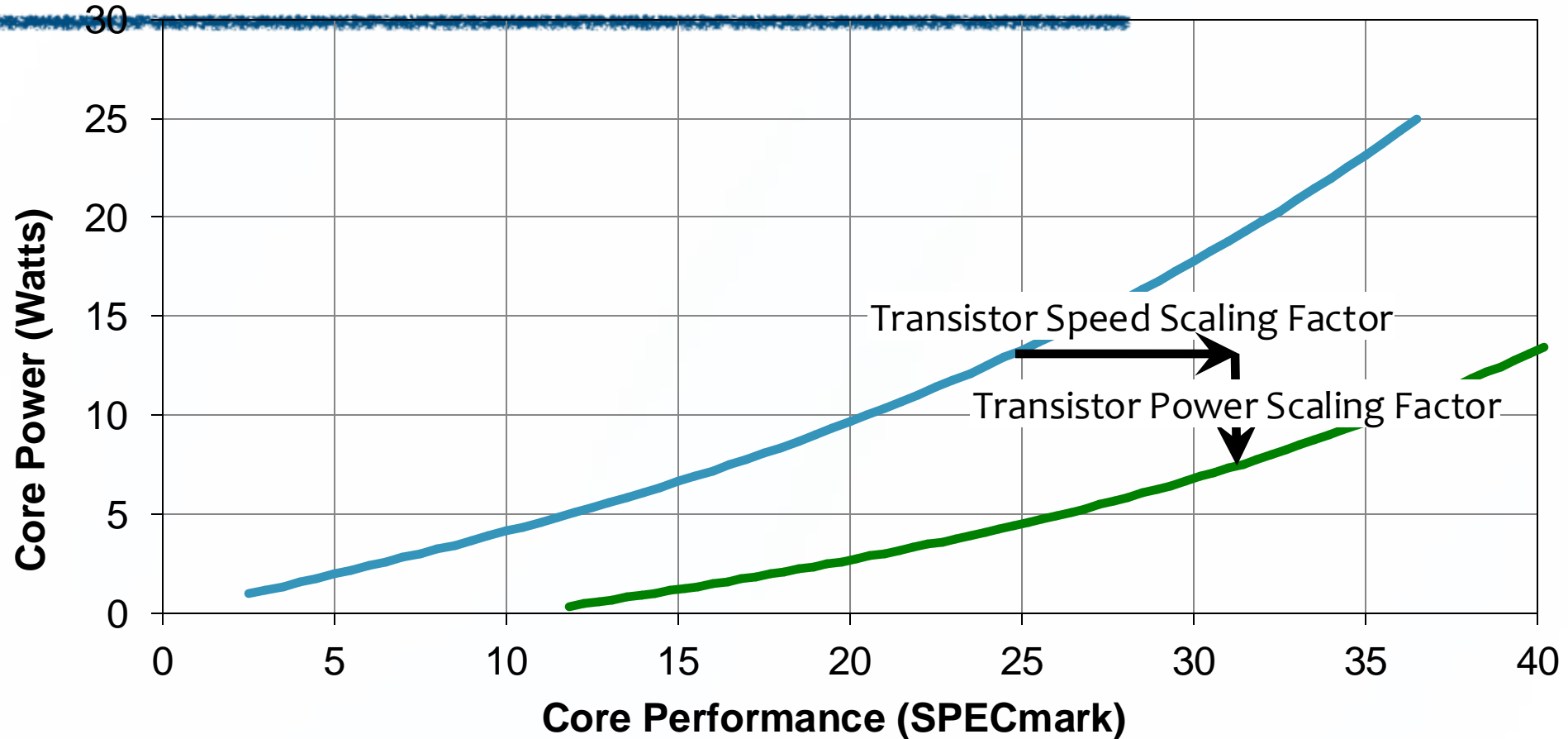
Single-core model (45 nm)



Power-Performance and Area-Performance
Pareto Optimal Frontiers

Single-core scaling model

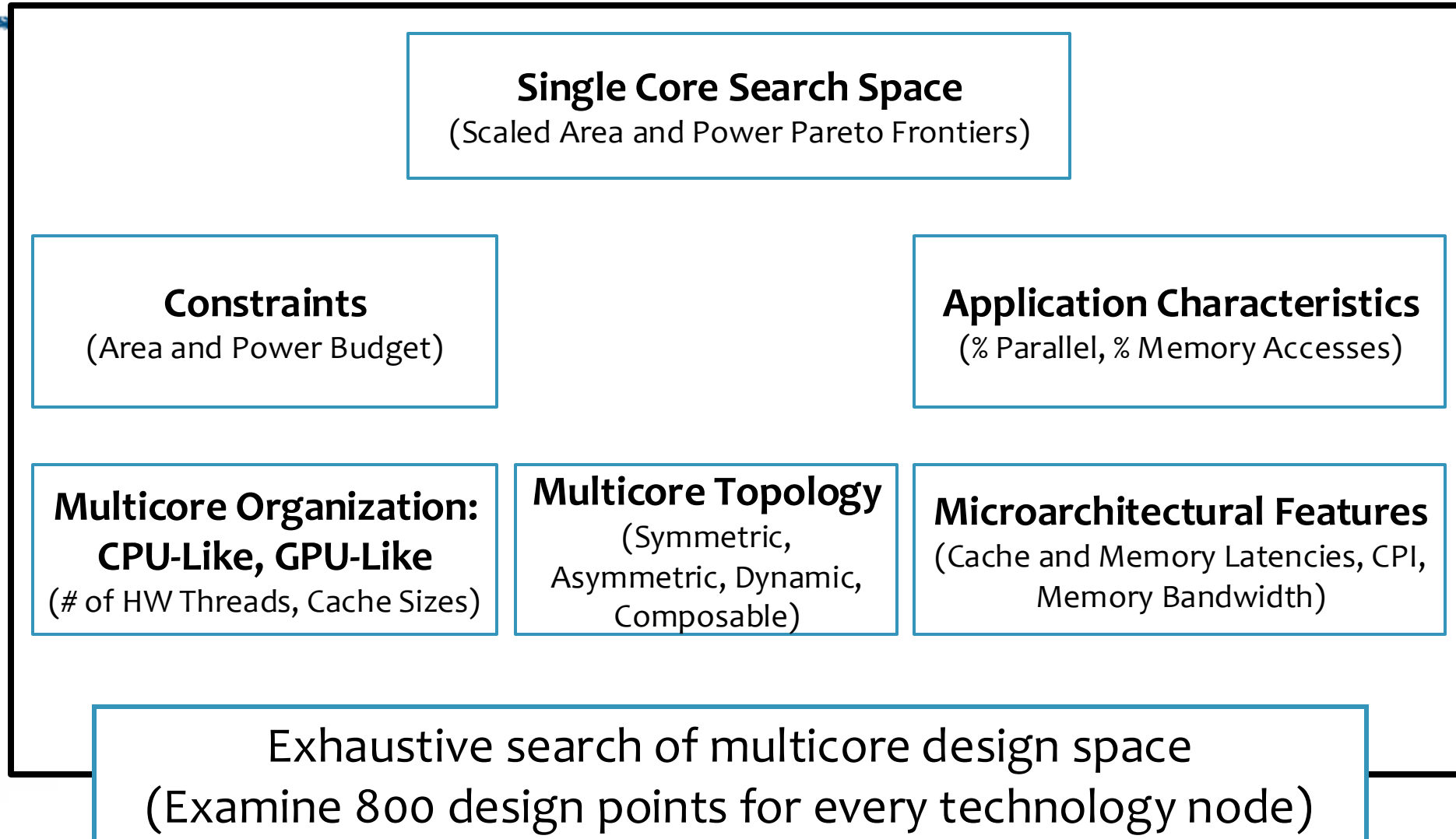
From 45 nm to 8 nm



Single-core Scaling Model:
Single-core Model \times Transistor Scaling Model

Multicore scaling model

From 45 nm to 8 nm



Multicore model (Amdahl's Law)

$$\text{Speedup} = \frac{1}{\frac{1 - f_{\text{Parallel}}}{\text{Serial Speedup}} + \frac{f_{\text{Parallel}}}{\text{Parallel Speedup}}}$$

Serial Speedup = $1 \times \text{Core Performance}$

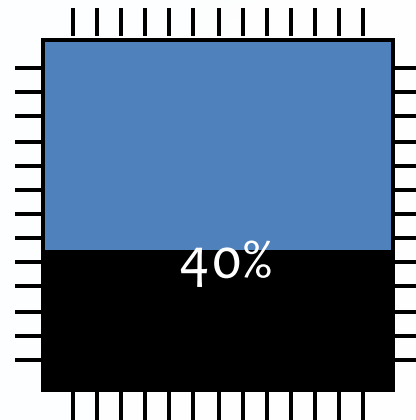
Parallel Speedup = $N \times \text{Core Performance}$



Dark silicon

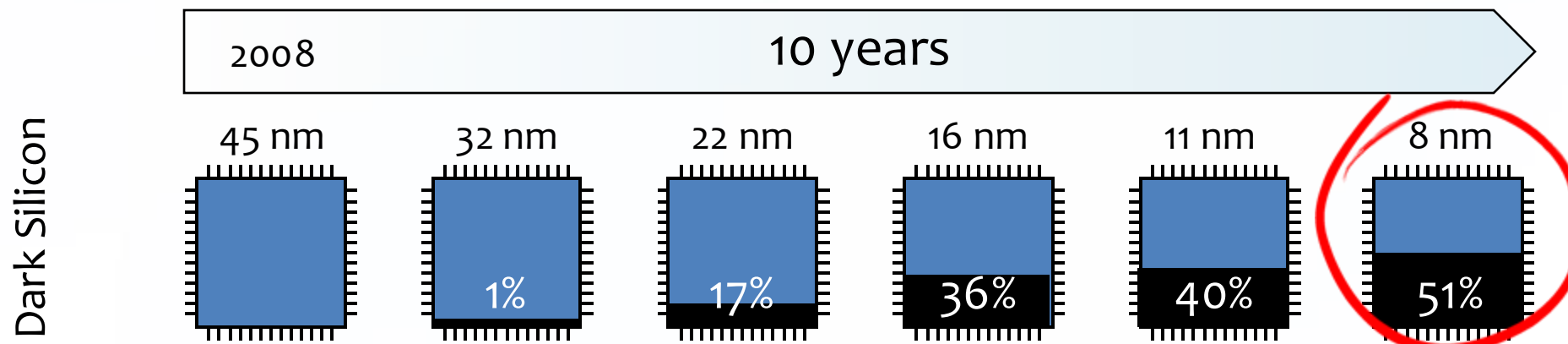
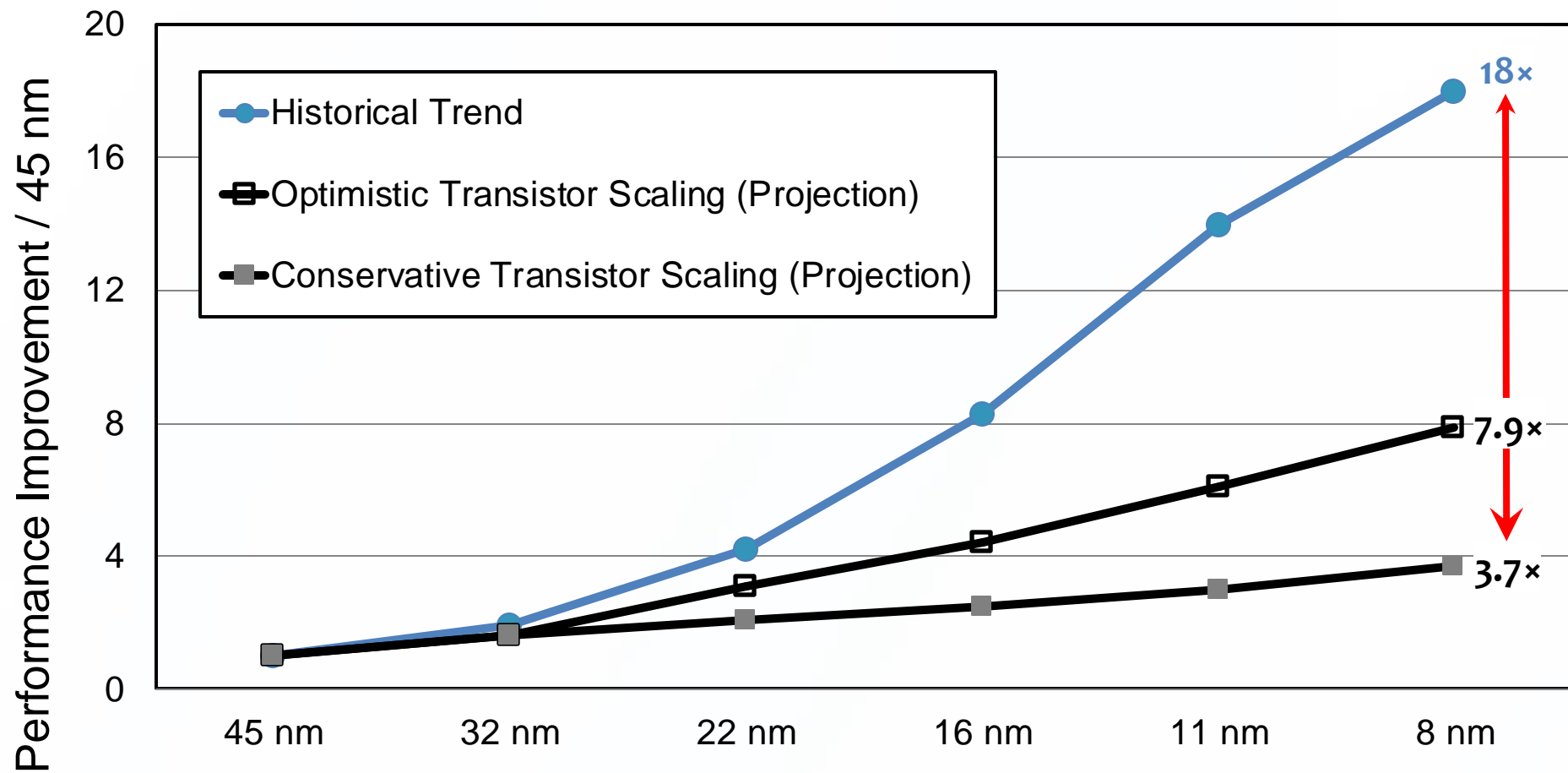
$$N_{Core} = \min\left(\frac{Area\ Budget}{Area_{Core}}, \frac{Power\ Budget}{Power_{Core}}\right)$$

$$Dark\ Silicon = 1 - \frac{N_{Core} \times Area_{Core}}{Area_{Budget}}$$



Evaluation Setup

- Applications:
 - 12 PARSEC Parallel Benchmarks
- Baseline:
 - The best multicore design available at 45 nm
- Constraints:
 - Driven from the best multicore design at 45 nm
 - Fixed Power Budget: 125 W
 - Fixed Area Budget: 111 mm²



The Shift Towards Domains-Specific Accelerators

Esmailzadeh et al. “Dark Silicon and the End of Multi-Core Scaling,” ISCA 2011

CACM Research Highlight

IEEE Micro Top Picks

The New York Times

Progress Hits Snag: Tiny Chips Use Outsize Power



Share full article



78

By John Markoff

July 31, 2011

For decades, the power of computers has grown at a staggering rate as designers have managed to squeeze ever more and ever tinier transistors onto a silicon chip — doubling the number every two years, on average, and leading the way to increasingly powerful and inexpensive personal computers, laptops and smartphones.

Industry of replacement?

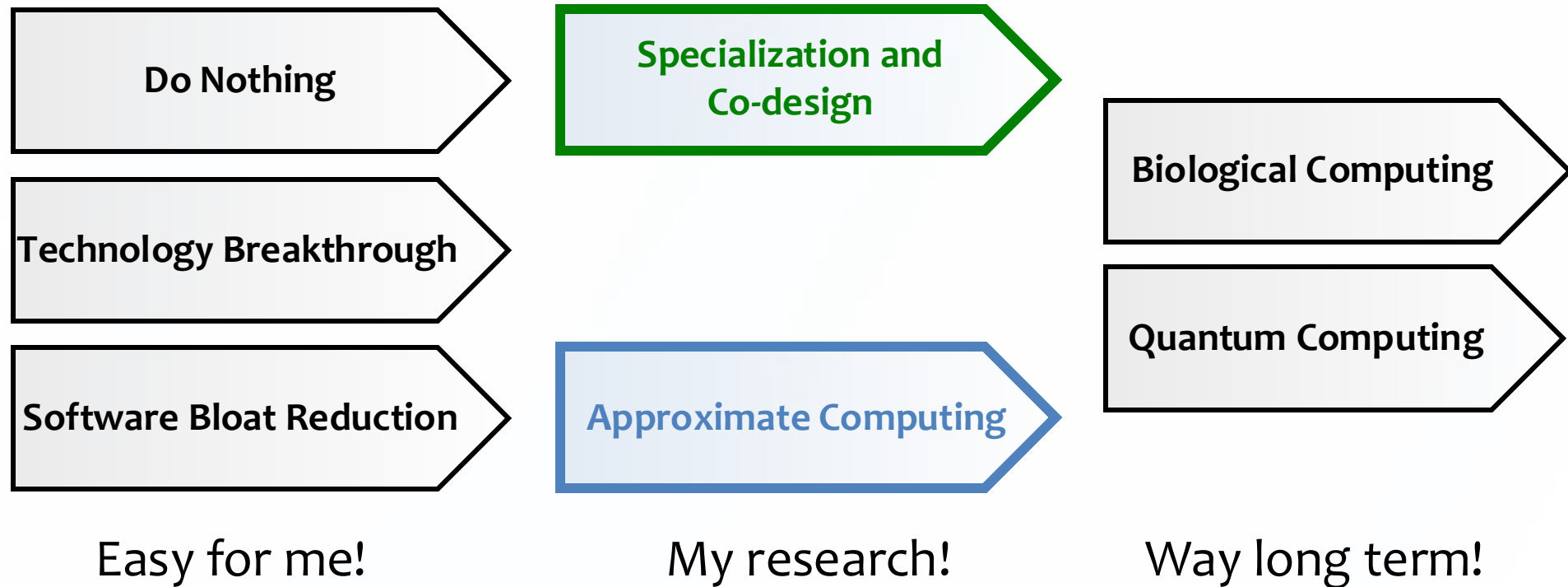
- Multicores are likely to be a stopgap
 - Not likely to continue the historical trends
 - Do not overcome the transistor scaling trends
 - The performance gap is significantly large
- Radical departures from conventional approaches are necessary
 - Extract more performance and efficiency from silicon while preserving programmability
 - Explore other sources of computing

Agenda

1. Who is Hadi
2. Course organization
3. Why CSE 240C Advanced Microarchitecture
 1. How we became an industry of new capabilities
 2. Why we might become an industry of replacement

3. Specialization and accelerators

Possible paths forward



Approximate computing

Embracing error

- Relax the abstraction of near-perfect accuracy in general-purpose computing
- Allow errors to happen in the computation
 - Run faster
 - Run more efficiently



WEB IMAGES VIDEOS MAPS NEWS MORE

bing

Hadi



Size ▾ Color ▾ Type ▾ Layout ▾ People ▾



New landscape of computing

Personalized and targeted computing

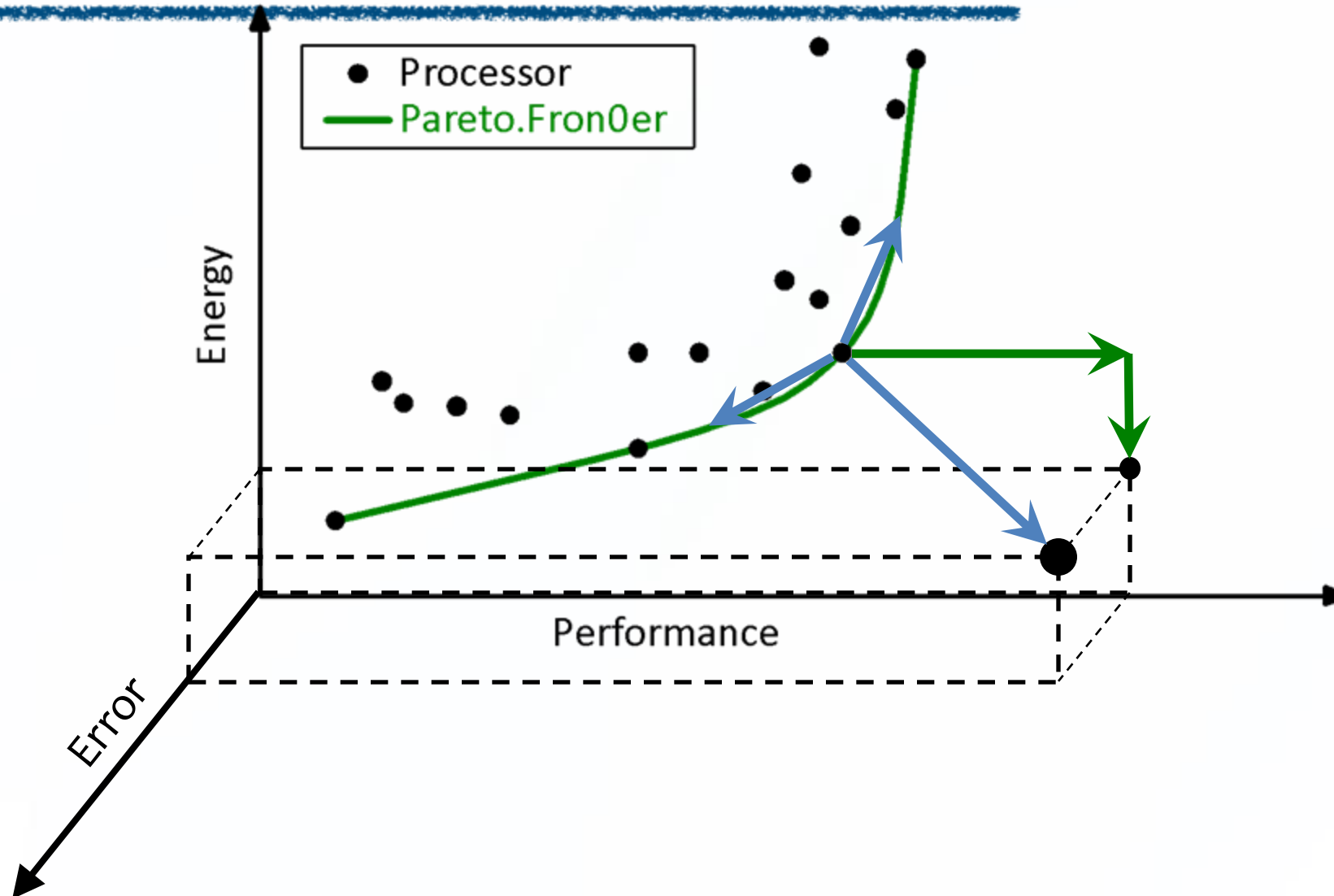


Classes of approximate applications

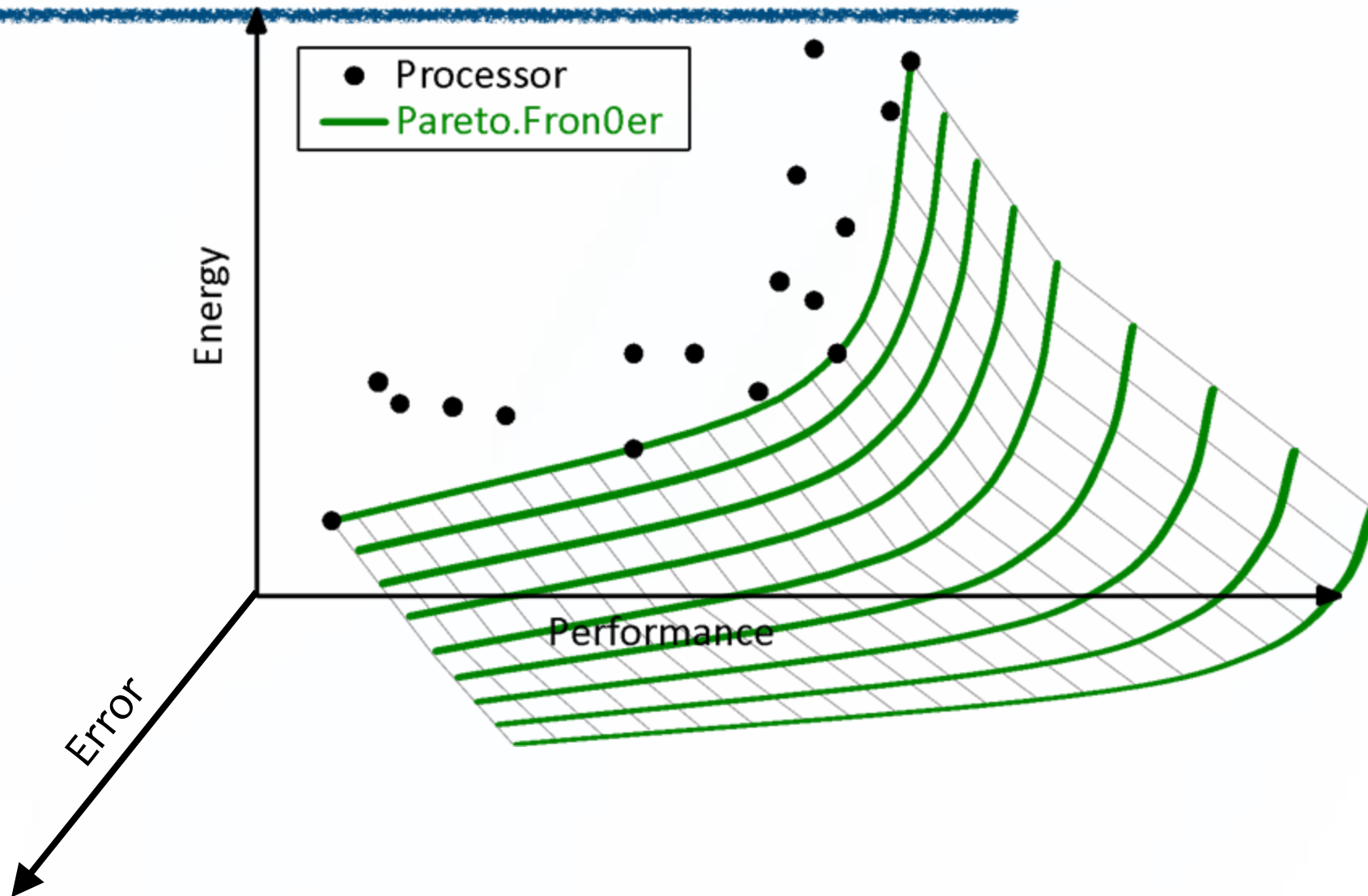
- Programs with analog inputs
 - Sensors, scene reconstruction
- Programs with analog outputs
 - Multimedia
- Programs with multiple possible answers
 - Web search, machine learning
- Convergent programs
 - Gradient descent, big data analytics

Adding a third dimension

Embracing Error



A fertile ground for innovation



Approximate computing techniques

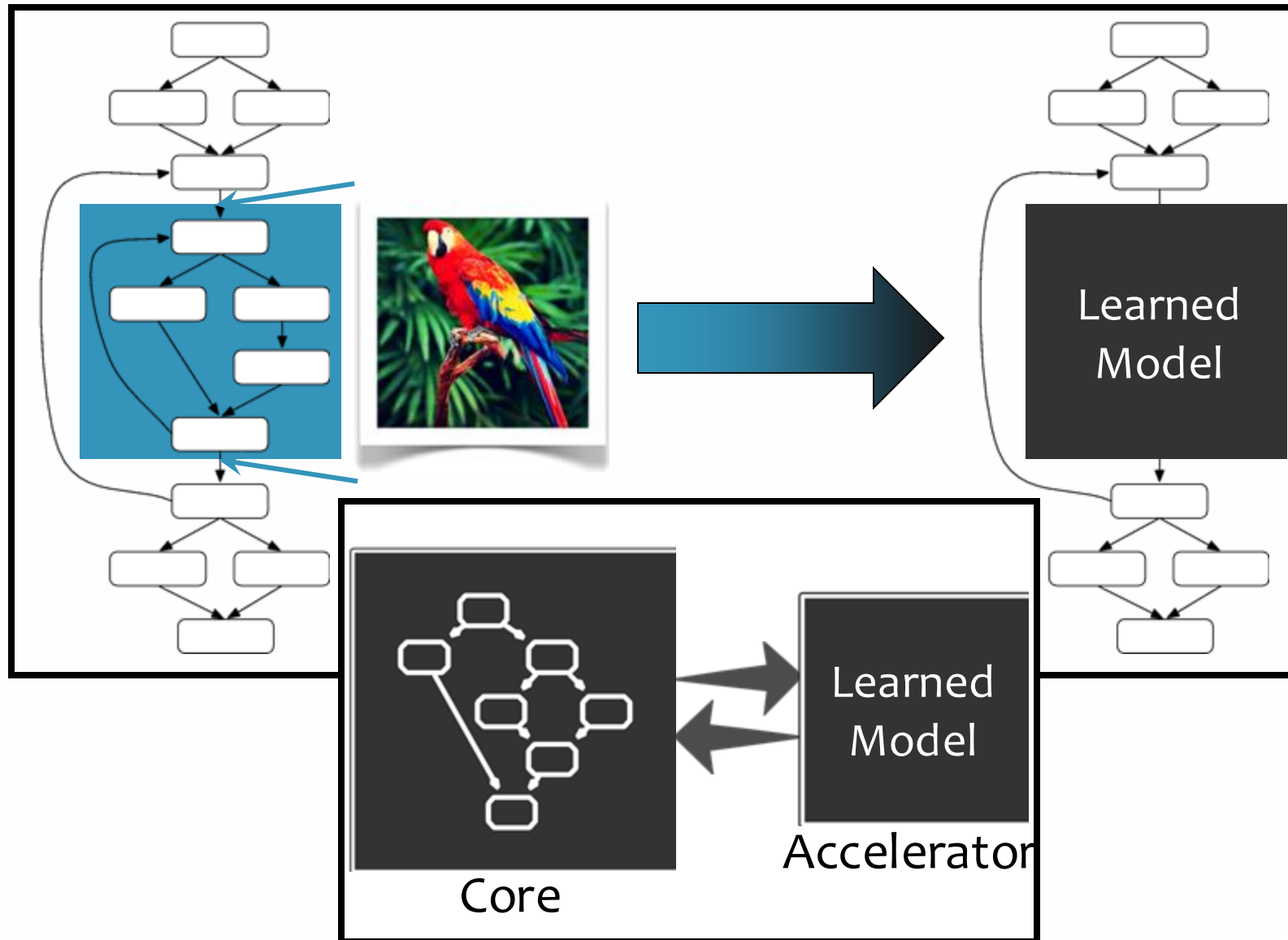
Same Model

- Sampling
 - Loop perforation (MIT)
- Compression
 - Sage (Michigan)
- Early termination
 - Green (MSR)
- Replacement
 - Green (MSR)
- Lower voltage
 - Truffle (Rice, UW)

From Model to Model

- von Neumann to Neural
 - NPUs (UW, GaTech, UCSD)

Parrot Algorithmic Transformation



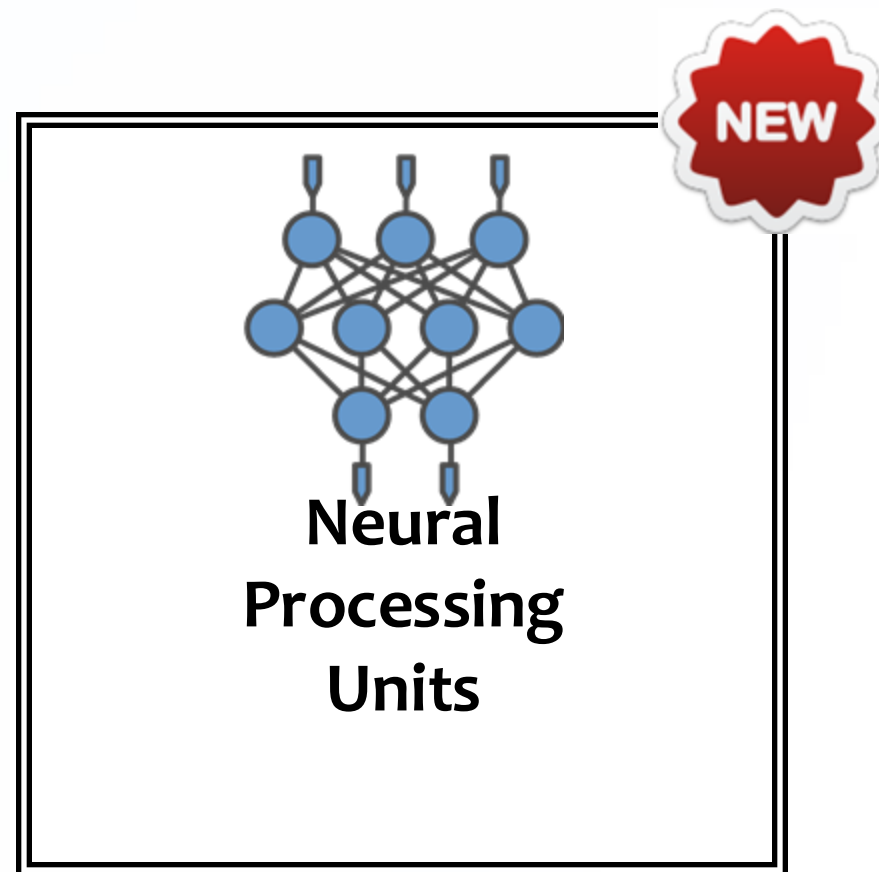
Neural Networks for Code Approximation

Powerful prediction tools

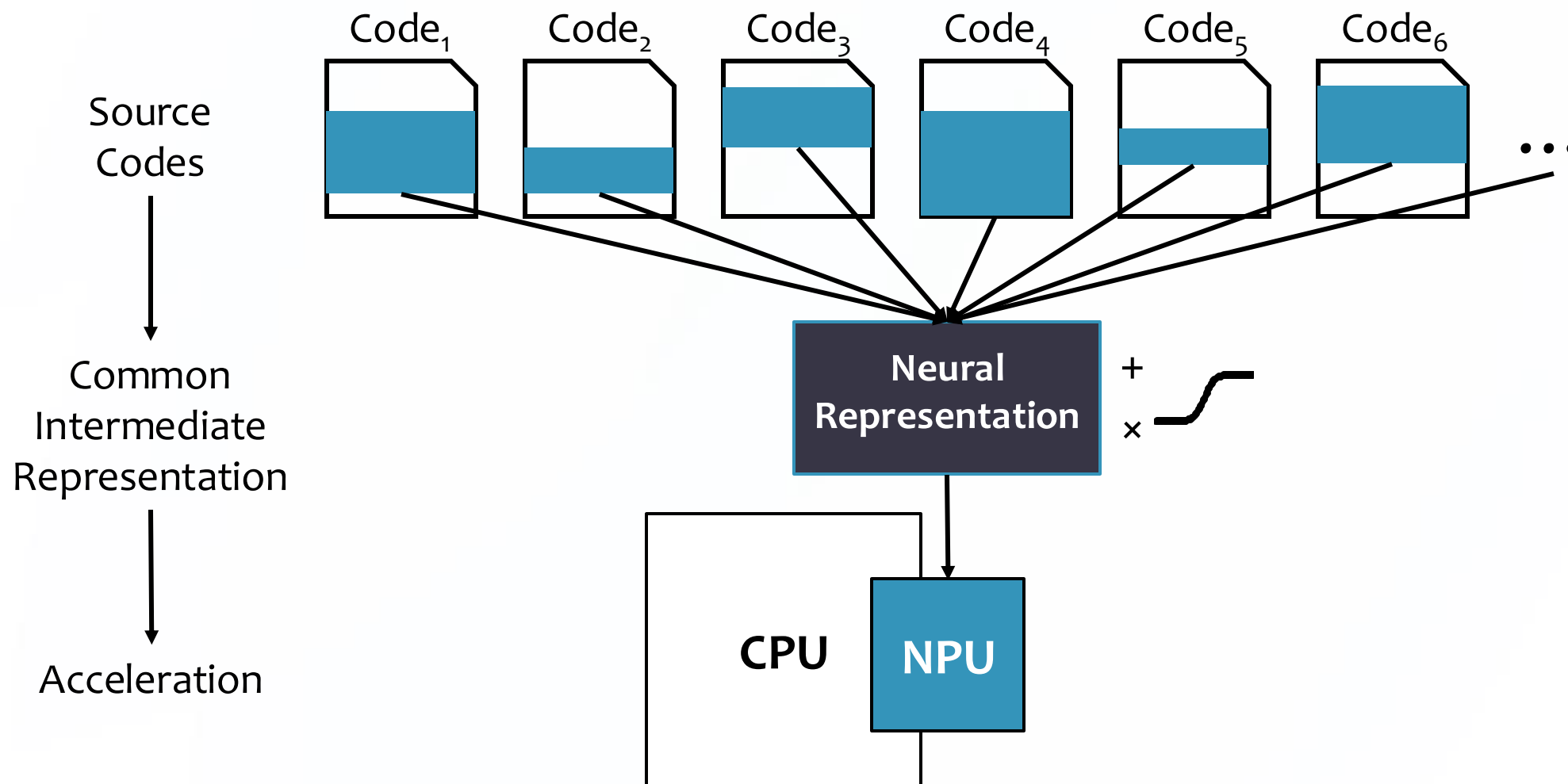
Highly parallel

Efficiently implementable with both
digital and analog hardware

Fault tolerant



NPU Acceleration

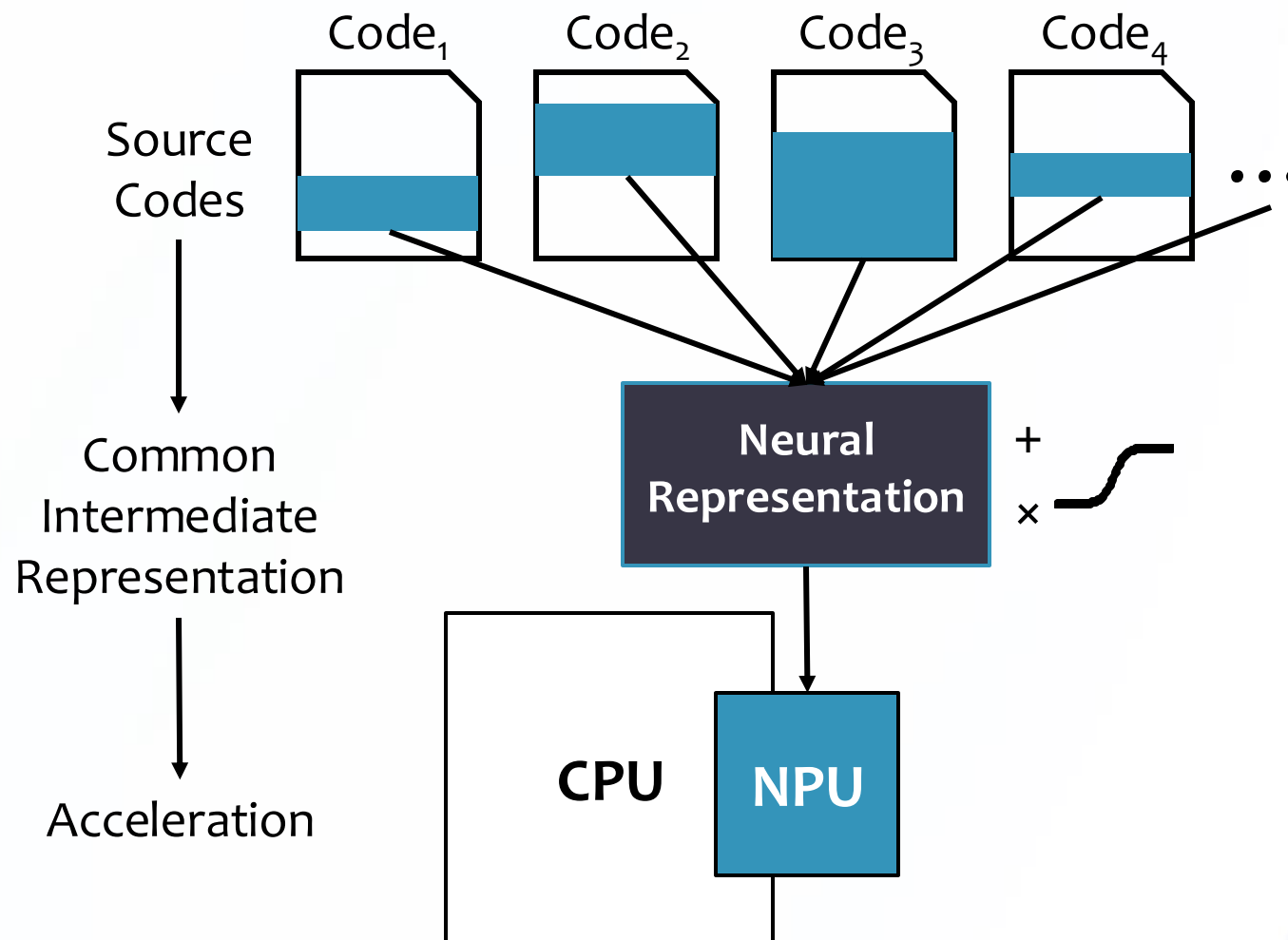


Neural Processing Units (NPUs)

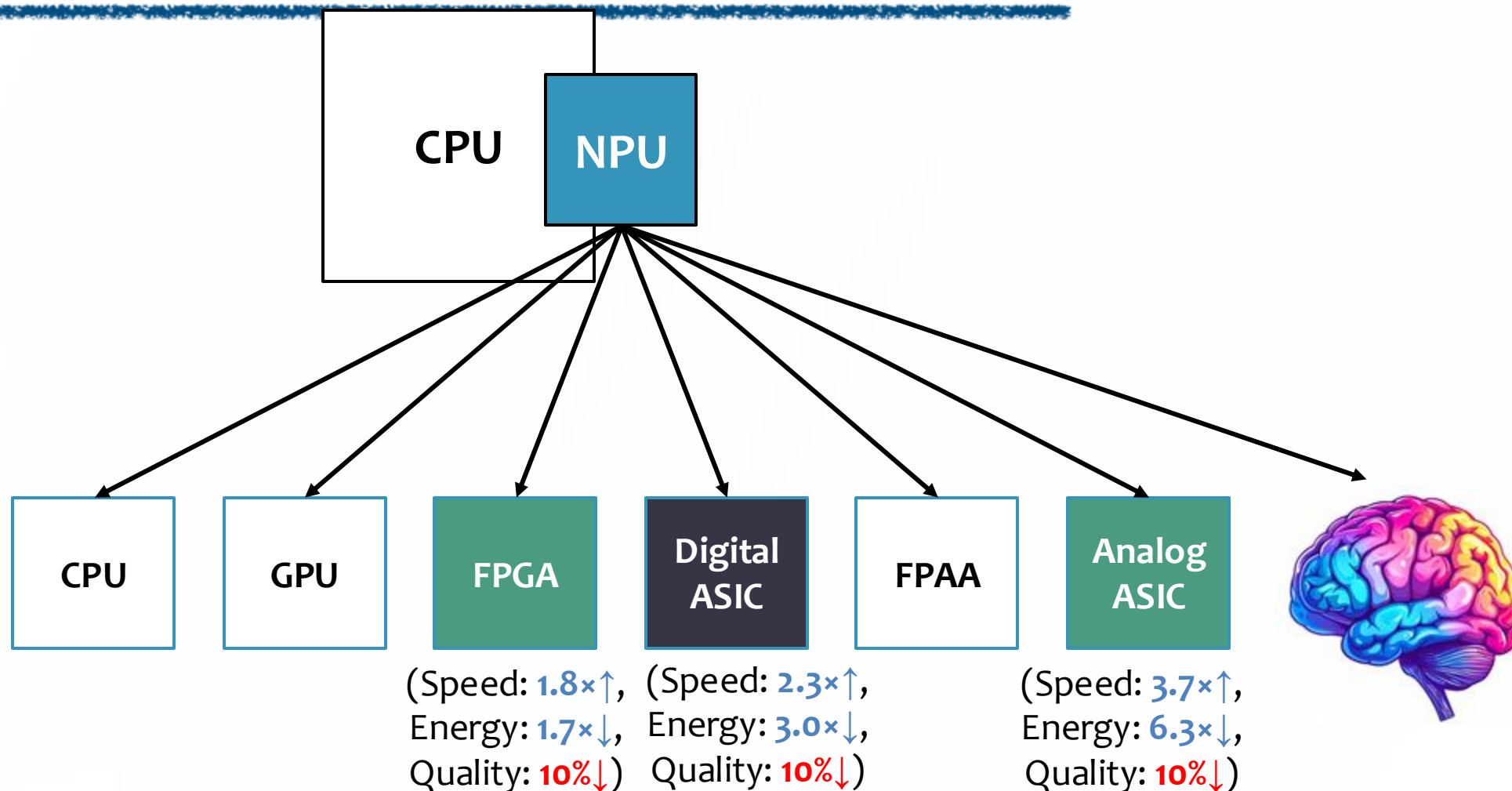
Esmailzadeh et al. "Neural Acceleration for General-Purpose Approximate Programs," Micro 2012

CACM Research Highlights

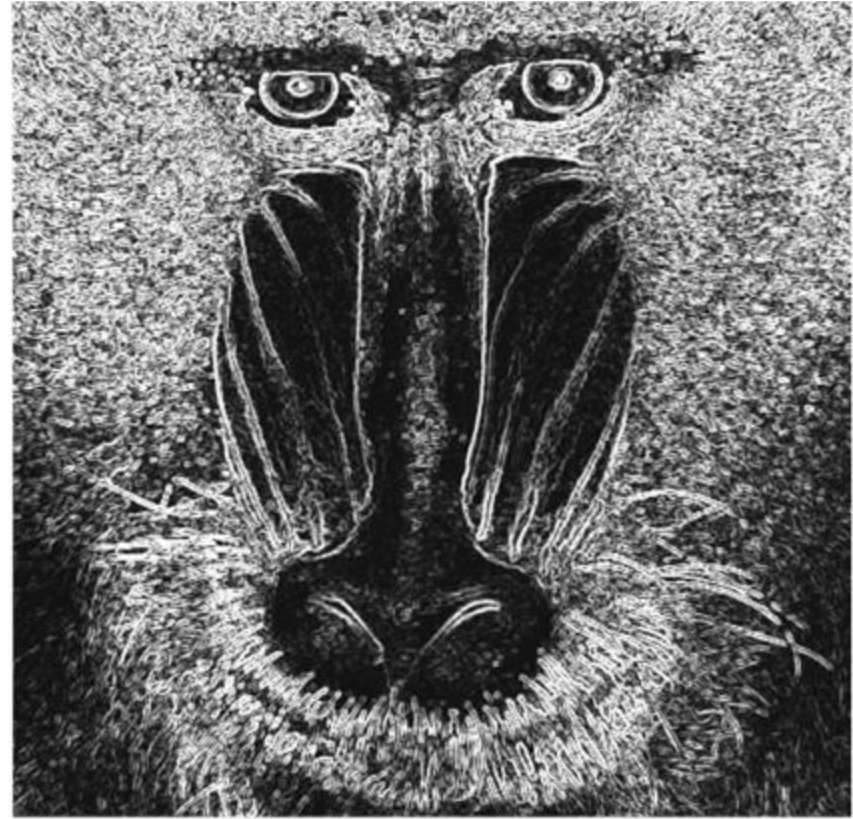
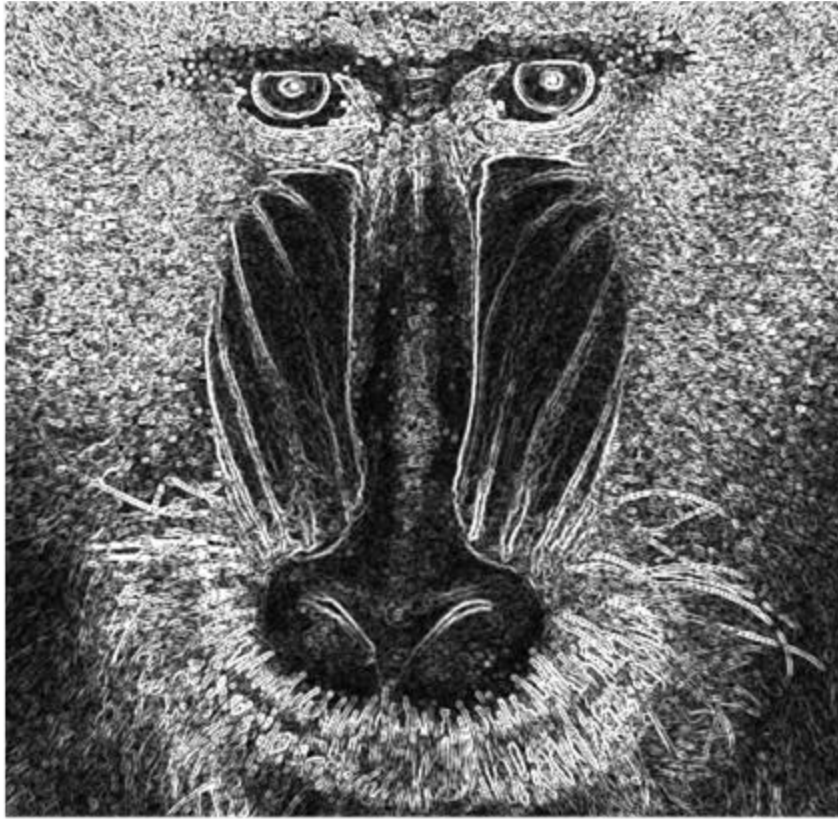
IEEE Micro Top Picks



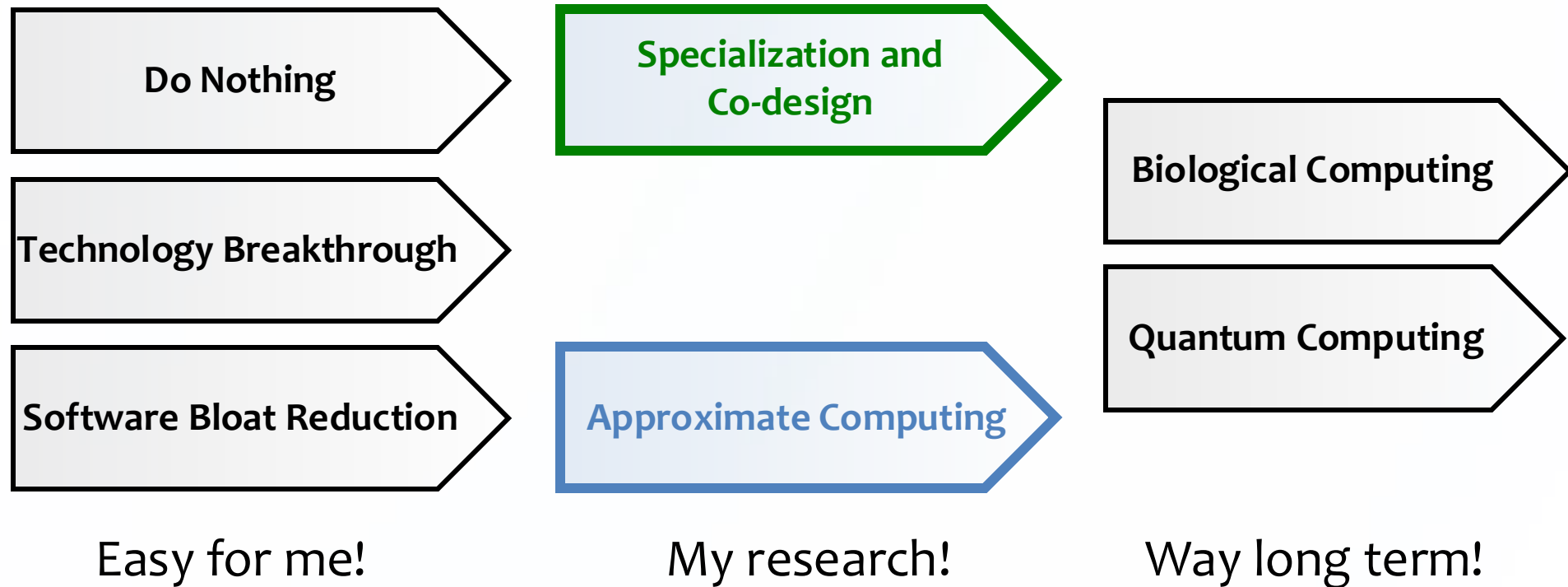
NPU design alternatives



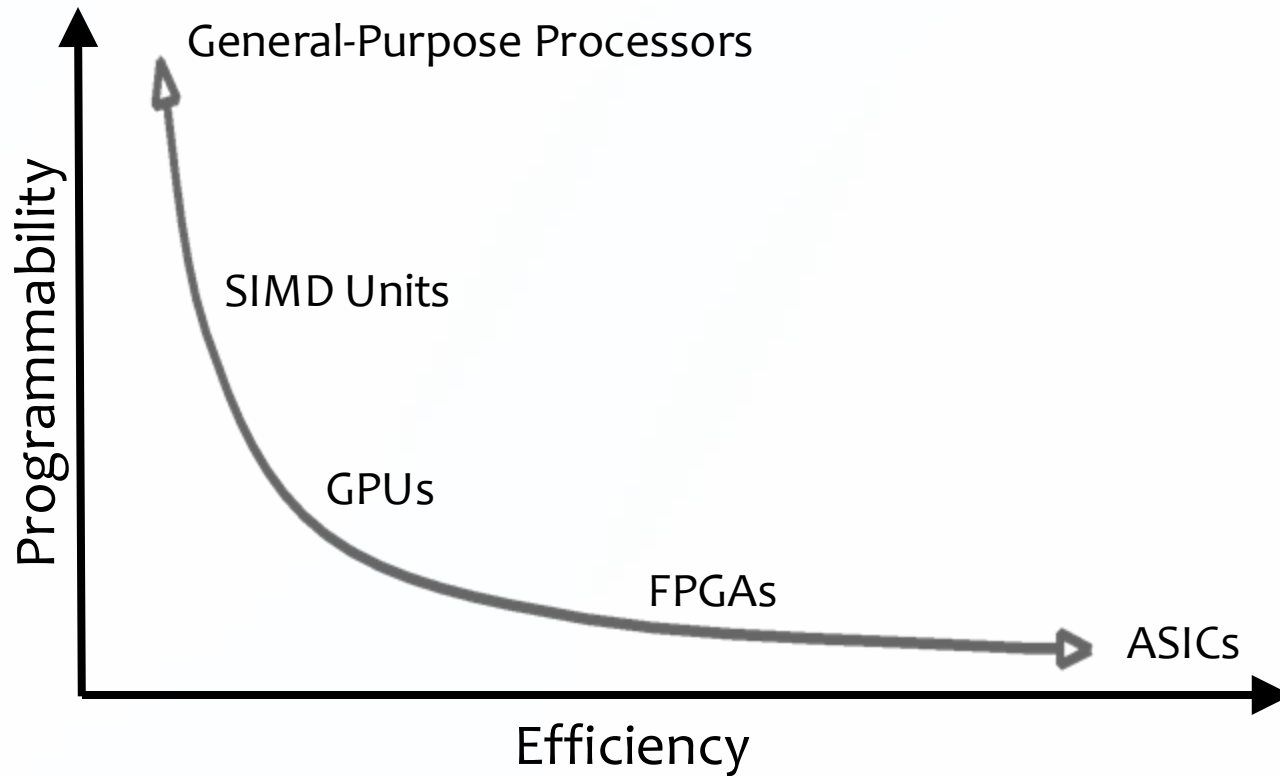
Approximate Computing versus Conventional Computing



Possible paths forward



Programmability versus Efficiency





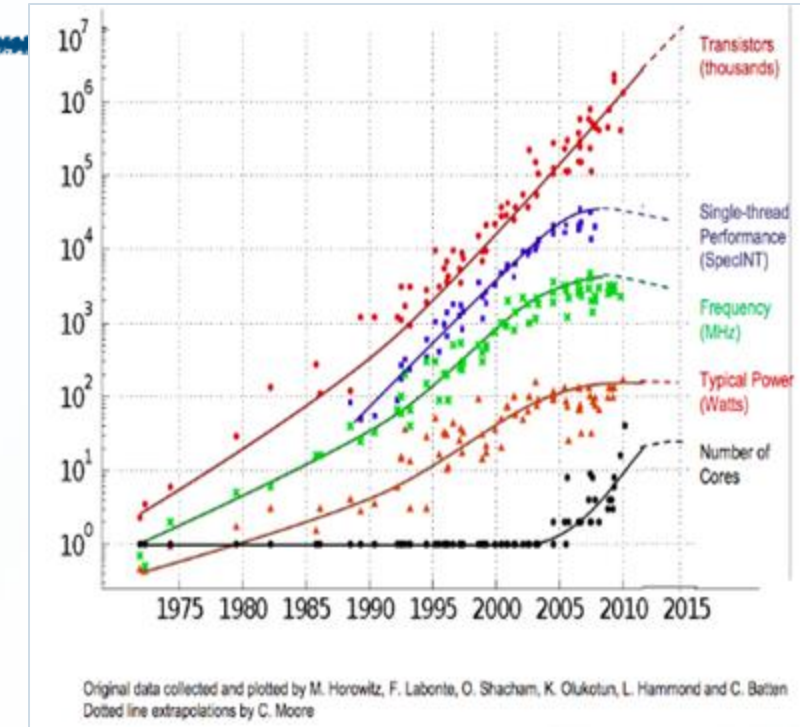
Large-Scale Reconfigurable Computing in a Microsoft
Datacenter

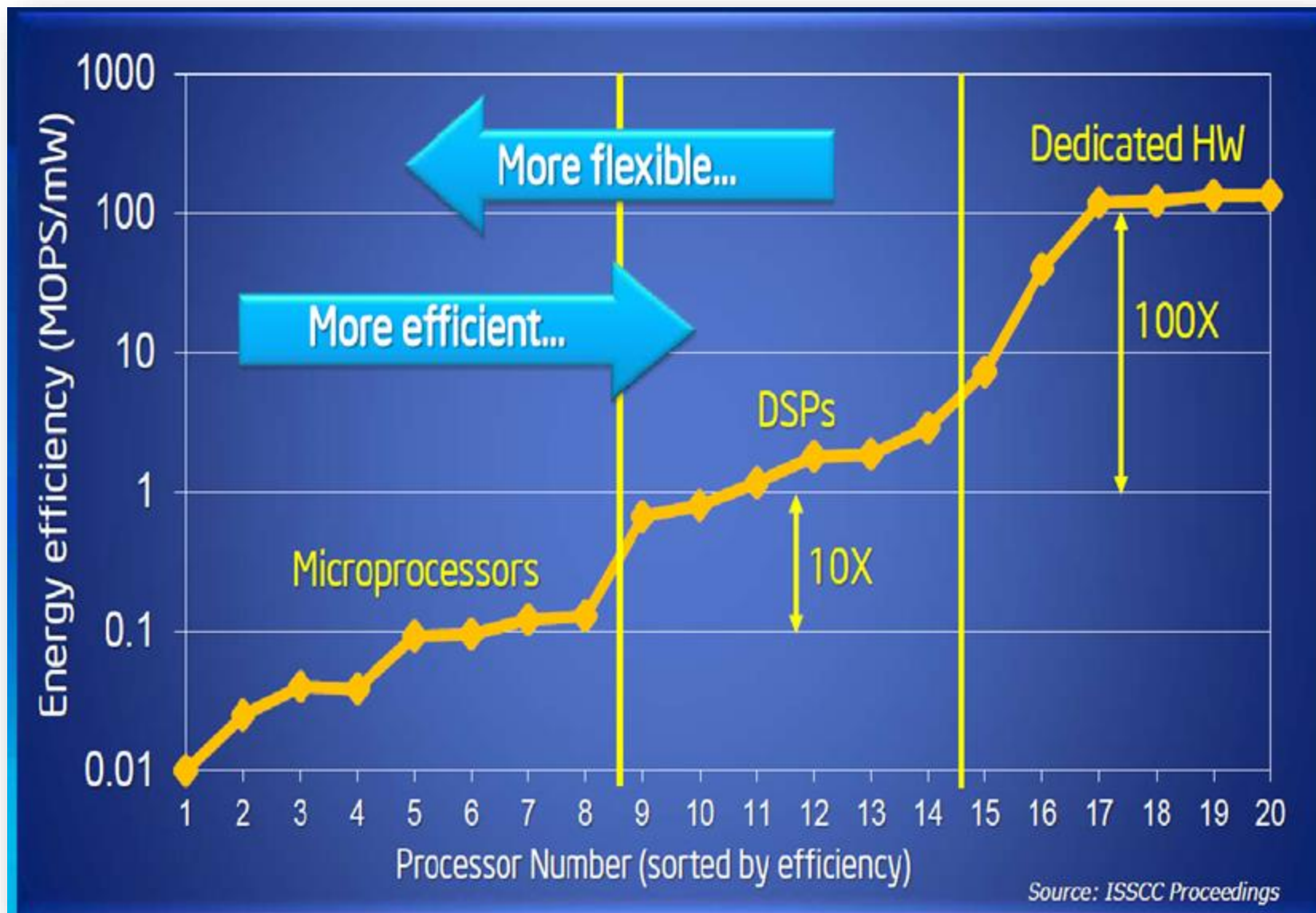


Microsoft Cloud Services



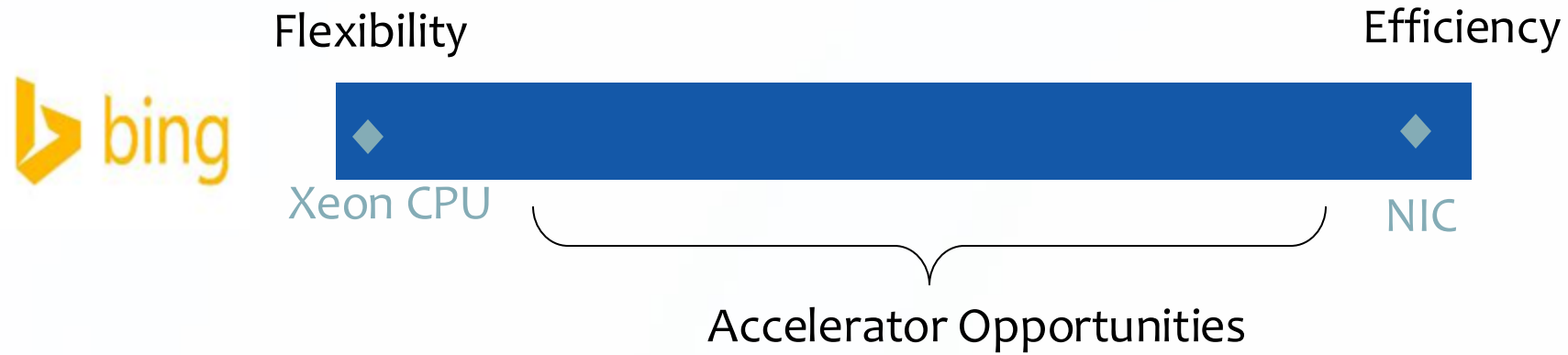
Capabilities, Costs



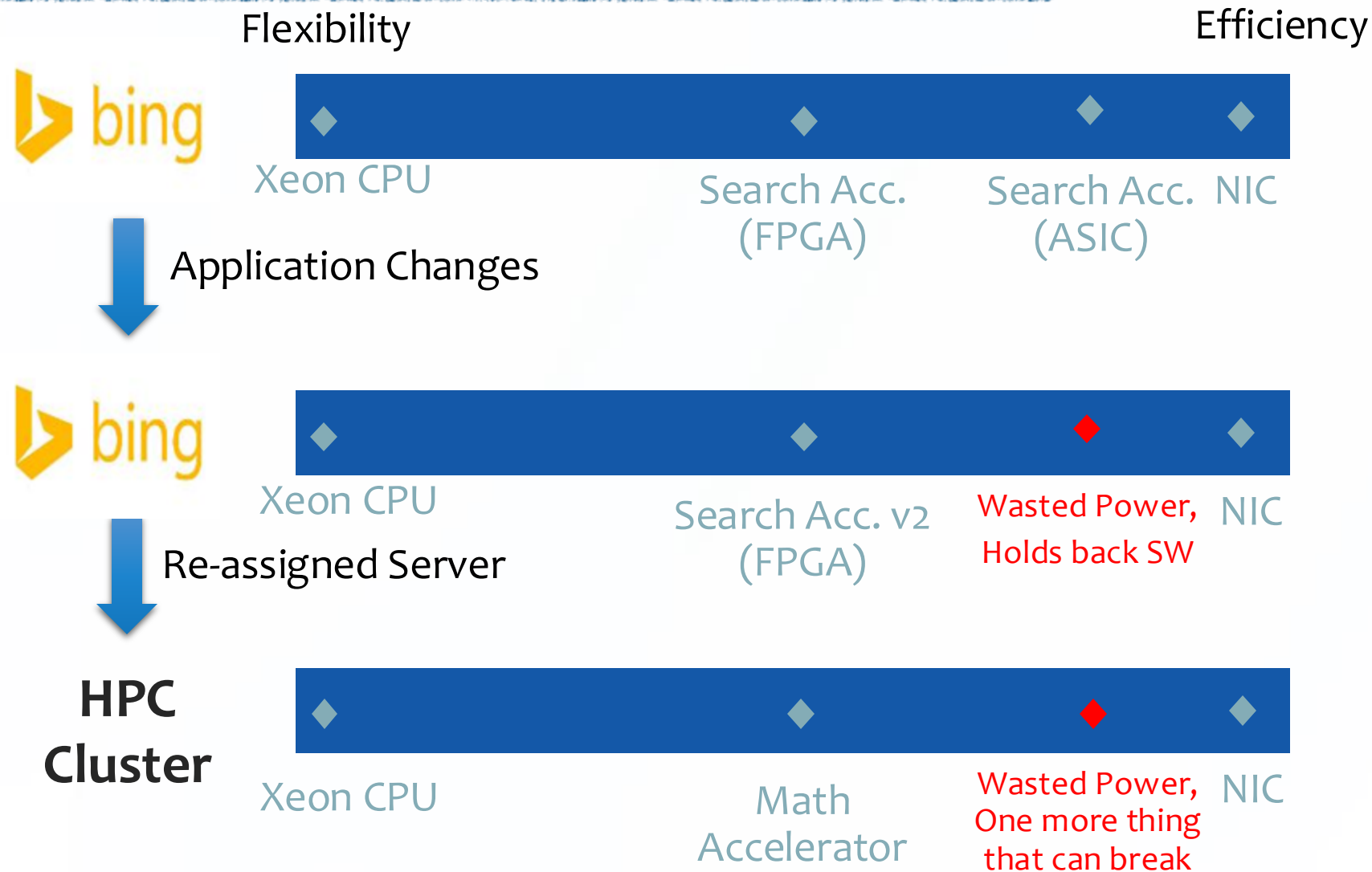


Increase Efficiency with Hardware Specialization

One Application's Accelerator



One Application's Accelerator



Integrating FPGAs into the Datacenter



This looks easy – Plug into the network and go!

Centralized



Distributed

Microsoft Open Compute Server



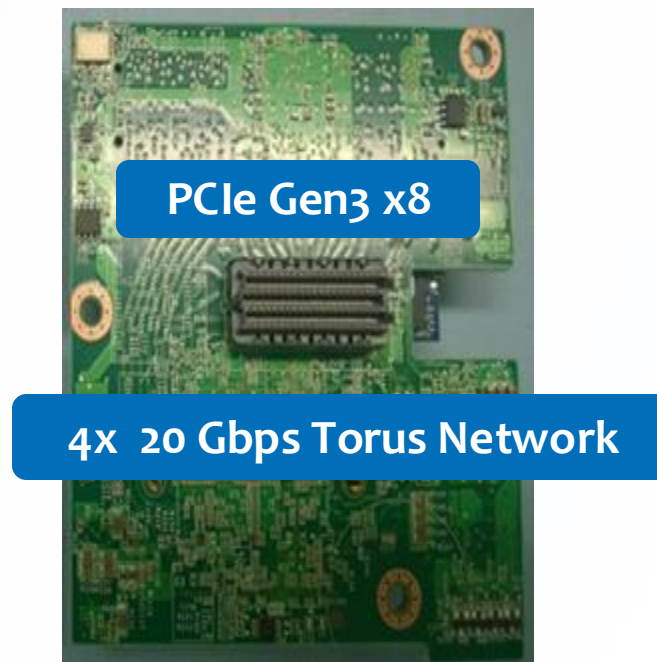
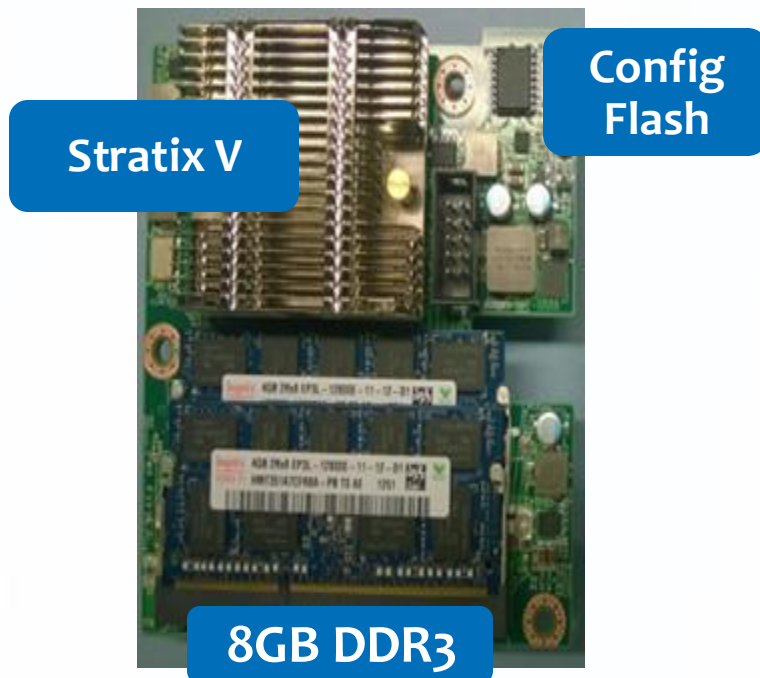
Two 8-core Xeon 2.1 GHz CPUs
64 GB DRAM
4 HDDs @ 2 TB, 2 SSDs @ 512 GB
10 Gb Ethernet
No cable attachments to server

Air flow

200 LFM
68 °C Inlet

Catapult FPGA Accelerator Card

- Altera Stratix V GS D5
 - 172k ALMs, 2,014 M20Ks, 1,590 DSPs
- 8GB DDR3-1333
- 32 MB Configuration Flash
- PCIe Gen 3 x8
- 8 lanes to Mini-SAS SFF-8088 connectors
- Powered by PCIe slot



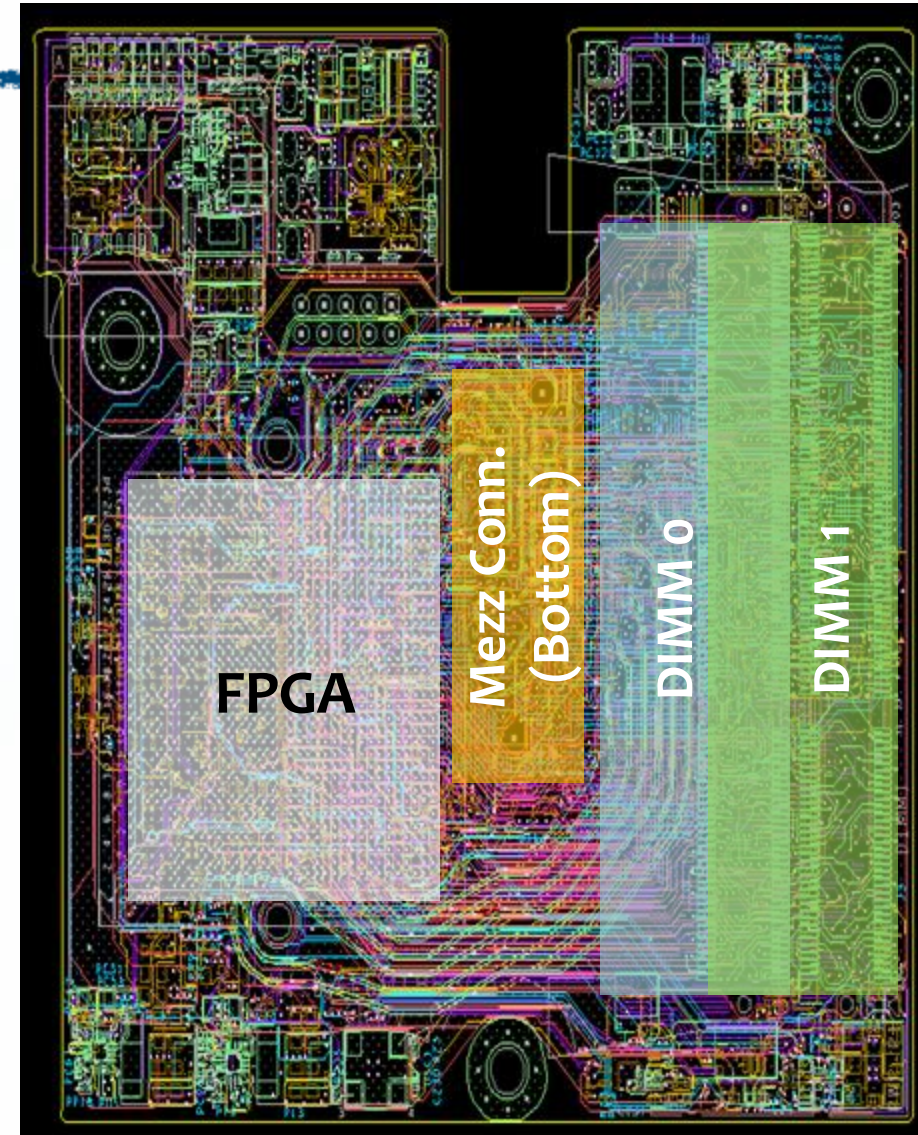
Board Details

16 Layer, FR408

9.5cm x 8.8cm x 115.8 mil

35mm x 35mm FPGA

14.2mm high heatsink



Scalable Reconfigurable Fabric

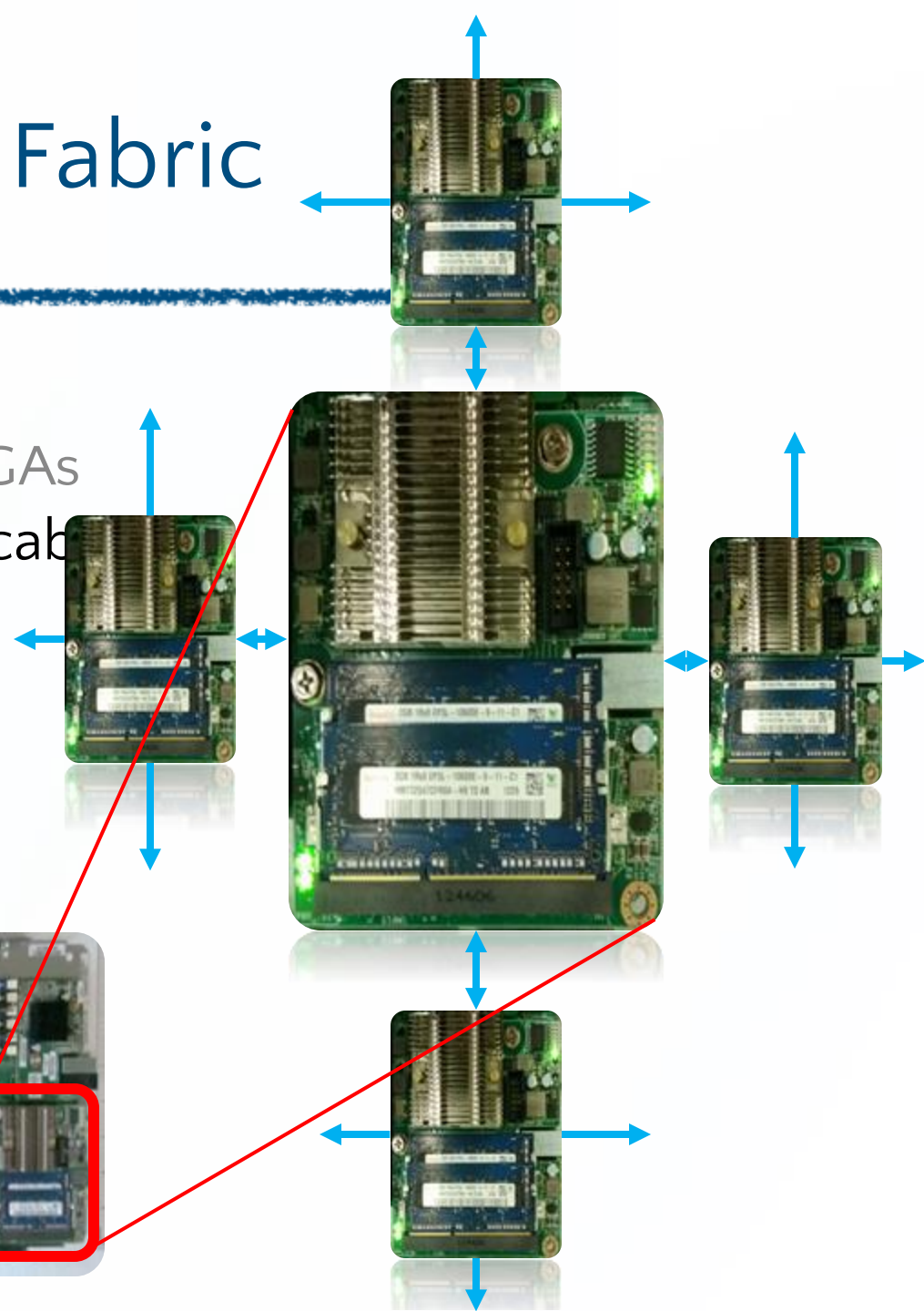
1 FPGA board per Server

48 Servers per ½ Rack

6x8 Torus Network among FPGAs

20 Gb over SAS SFF-8088 cable

Data Center Server (1U, ½ width)





Microsoft Research

@MSFTResearch



Follow

Catapult propels datacenter services into the future @Bing @MSFTResearch @dcburger #FPGA bit.ly/1lzp10f

Reply Retweet Favorite More



RETWEETS

56

FAVORITES

30



3:00 PM - 16 Jun 2014

Flag media

Reply to @MSFTResearch @bing @dcburger

Intel's \$16.7 Billion Altera Deal Is Fueled by Data Centers

by Ian King

June 1, 2015 — 8:34 AM EDT Updated on June 1, 2015 — 4:13 PM EDT



■ Intel Acquiring Altera in \$16.7B Chipmaker Combination

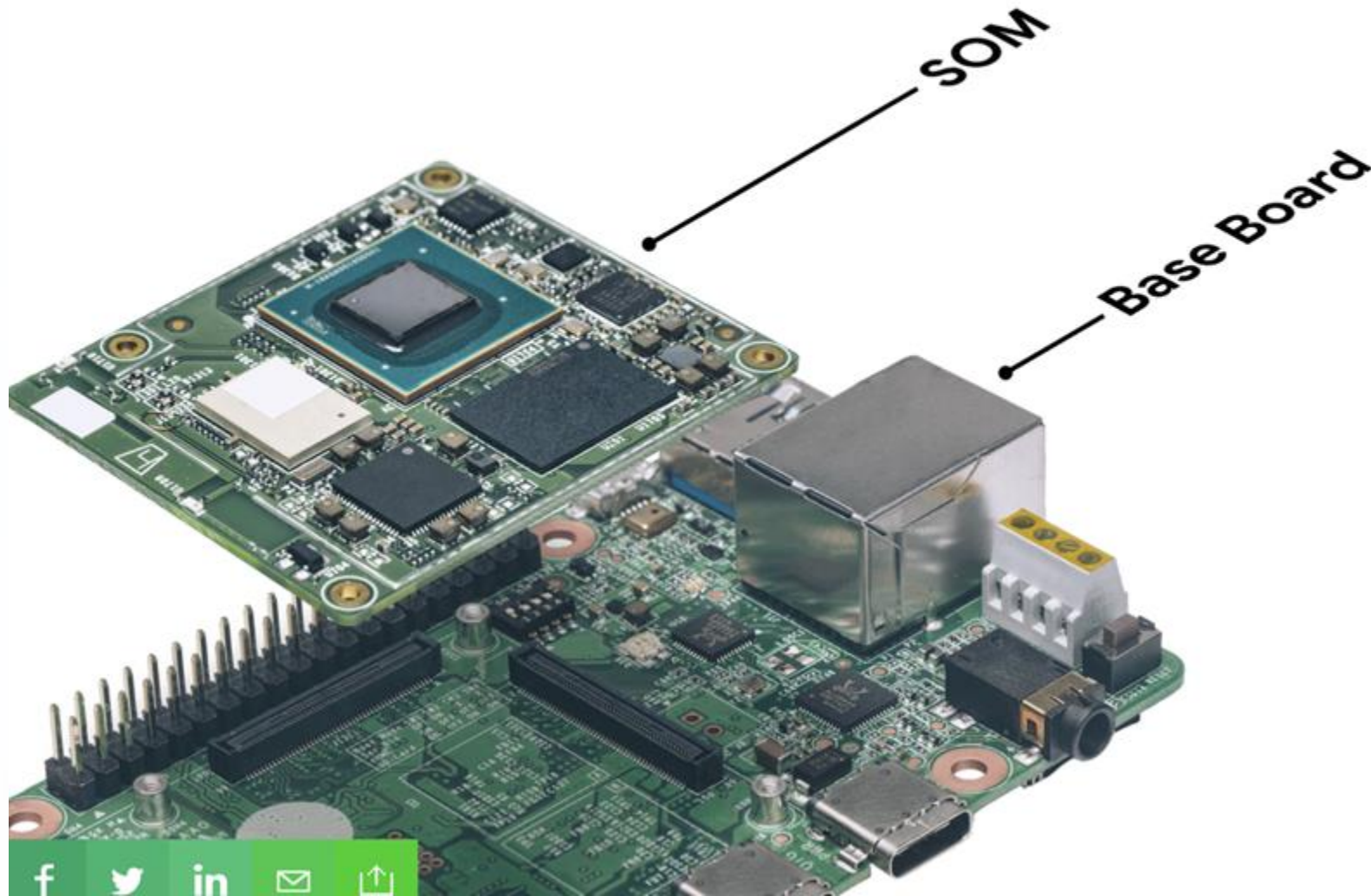


Intel Corp. agreed to buy Altera Corp. for \$16.7 billion to defend its presence in data centers, forging a deal that will add to a record year for industry consolidation.

Google is making a fast specialized TPU chip for edge devices and a suite of services to support it

Matthew Lynley @mattlynley / 2 months ago

 Comment

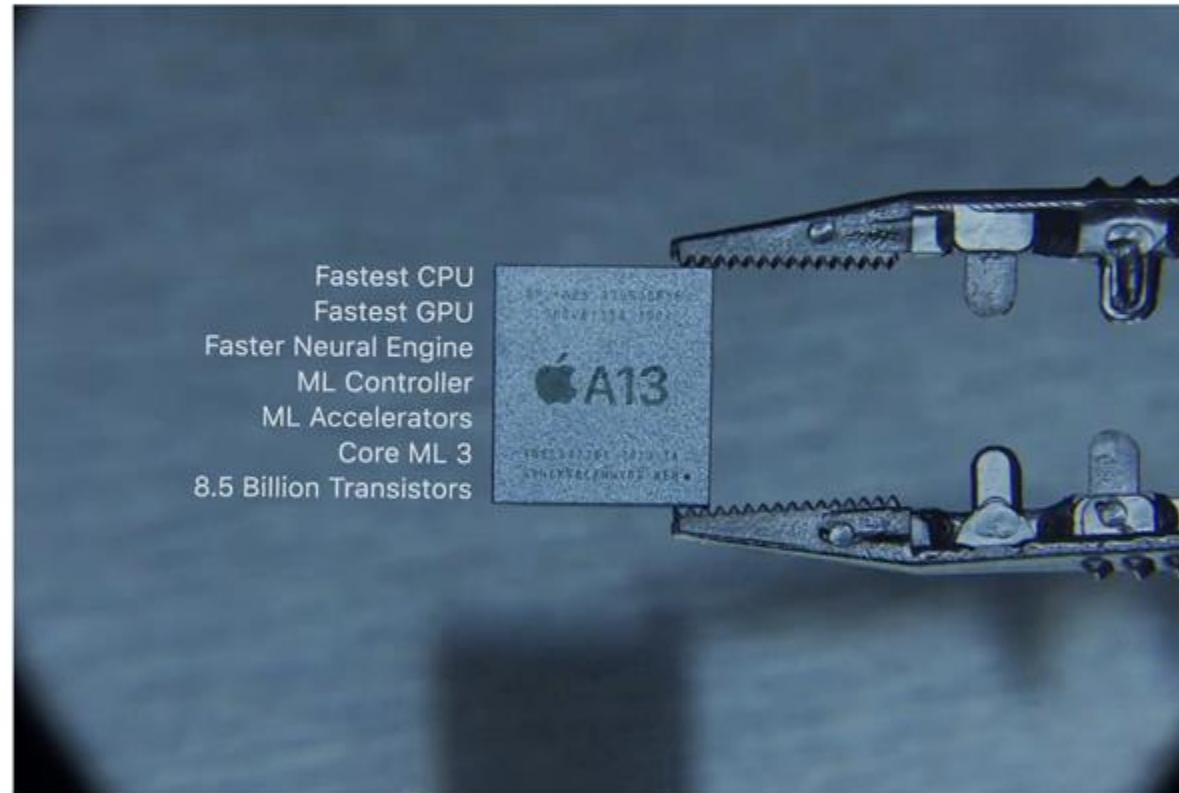


Apple says its new A13 Bionic chip brings hours of extra battery life to new iPhones

Plus a 20 percent performance boost across the board

By [Sean Hollister](#) | [@StarFire2258](#) | Sep 10, 2019, 2:02pm EDT

[f](#) [t](#) [SHARE](#)



Apple has revealed the chip that will power [its new 2019 iPhones](#): the A13 Bionic. And as you'd expect, the company is wasting no time in explaining that it's the most powerful silicon ever to grace the inside of a smartphone — just as it has every year for the past three years. But if you care about battery life, you'll want to pay attention.



Amazon
free at B



The Fifth Day of Creation ...

By Maestro Mahmoud Farshchian

