

# Seminar in Computer Architecture Meeting 1: Introduction

Prof. Onur Mutlu

ETH Zürich

Fall 2020

17 September 2020



# Brief Self Introduction

---



## ■ Onur Mutlu

- ❑ Full Professor @ ETH Zurich ITET (INFK), since September 2015
- ❑ Strecker Professor @ Carnegie Mellon University ECE/CS, 2009-2016, 2016-...
- ❑ PhD from UT-Austin, worked at Google, VMware, Microsoft Research, Intel, AMD
- ❑ <https://people.inf.ethz.ch/omutlu/>
- ❑ [omutlu@gmail.com](mailto:omutlu@gmail.com) (Best way to reach me)
- ❑ <https://people.inf.ethz.ch/omutlu/projects.htm>

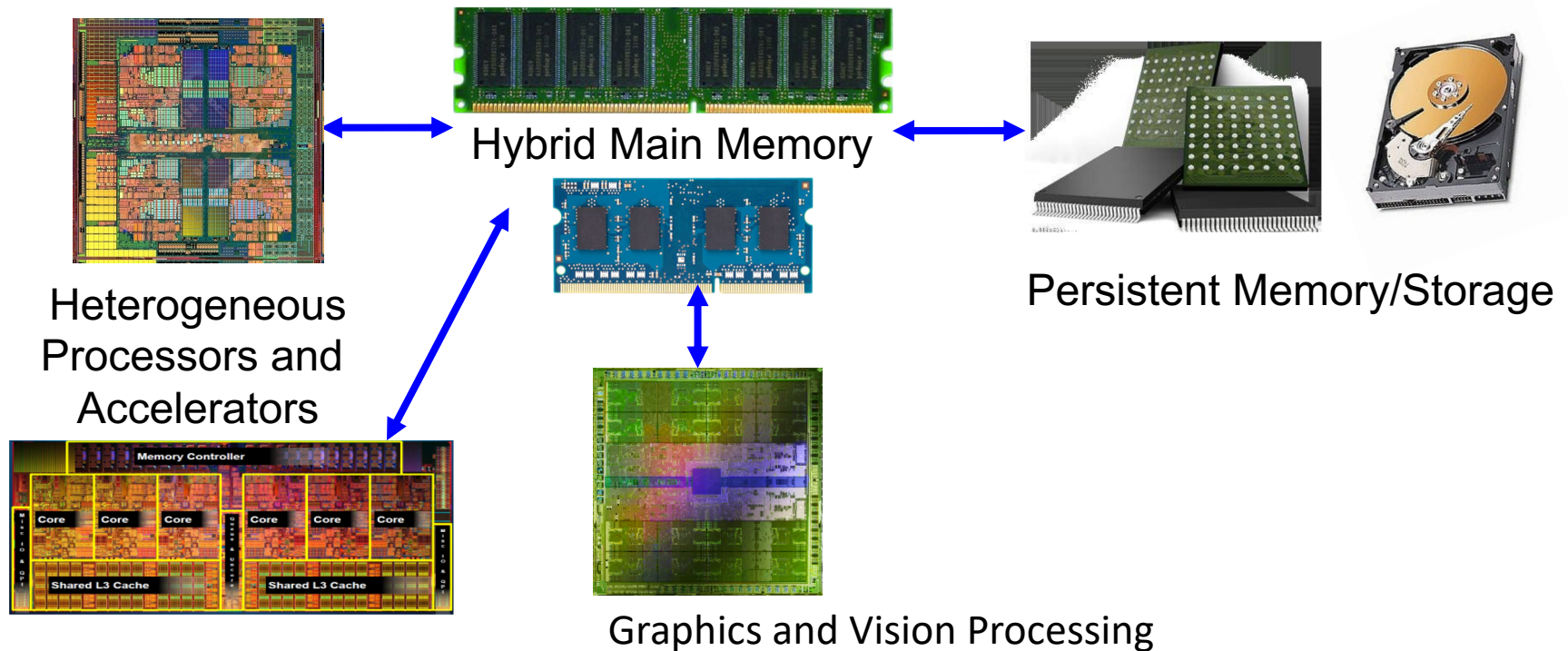
## ■ Research and Teaching in:

- ❑ Computer architecture, computer systems, hardware security, bioinformatics
- ❑ Memory and storage systems
- ❑ Hardware security, safety, predictability
- ❑ Fault tolerance
- ❑ Hardware/software cooperation
- ❑ Architectures for bioinformatics, health, medicine
- ❑ ...

# Current Research Mission

---

*Computer architecture, HW/SW, systems, bioinformatics, security*



**Build fundamentally better architectures**

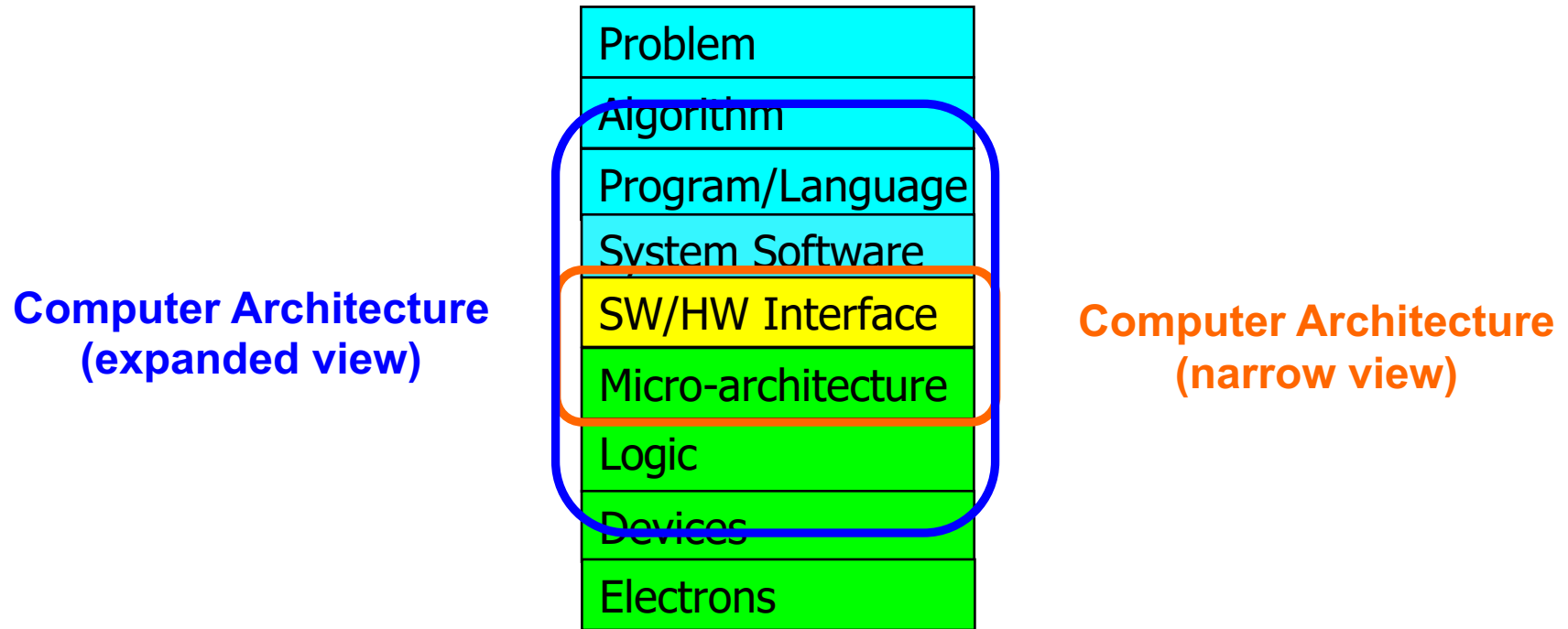
# Four Key Current Directions

---

- Fundamentally **Secure/Reliable/Safe** Architectures
- Fundamentally **Energy-Efficient** Architectures
  - **Memory-centric** (Data-centric) Architectures
- Fundamentally **Low-Latency and Predictable** Architectures
- Architectures for **AI/ML, Genomics, Medicine, Health**

# The Transformation Hierarchy

---

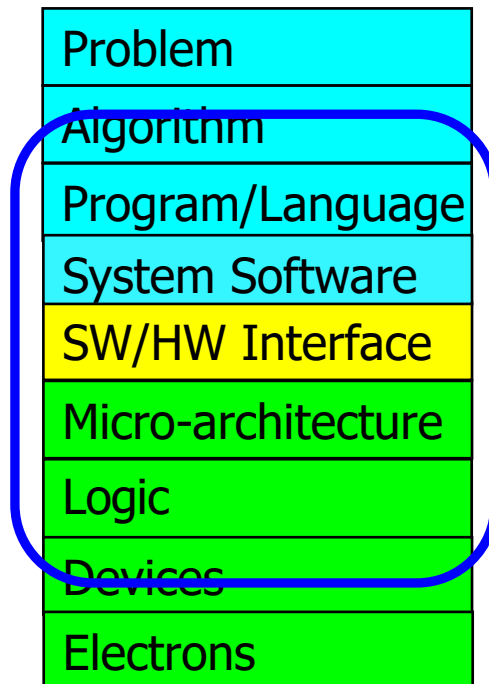


# Axiom

---

To achieve the highest **energy efficiency** and **performance**:

**we must take the expanded view**  
of computer architecture

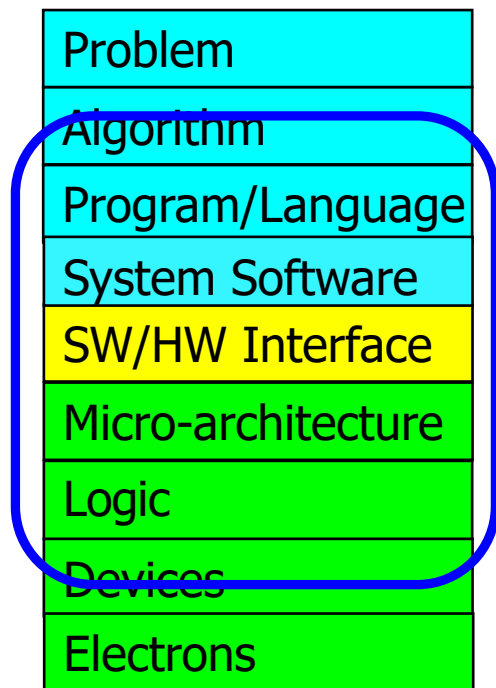


**Co-design across the hierarchy:**  
**Algorithms to devices**

**Specialize as much as possible**  
**within the design goals**

# Current Research Mission & Major Topics

## Build fundamentally better architectures



**Broad research  
spanning apps, systems, logic  
with architecture at the center**

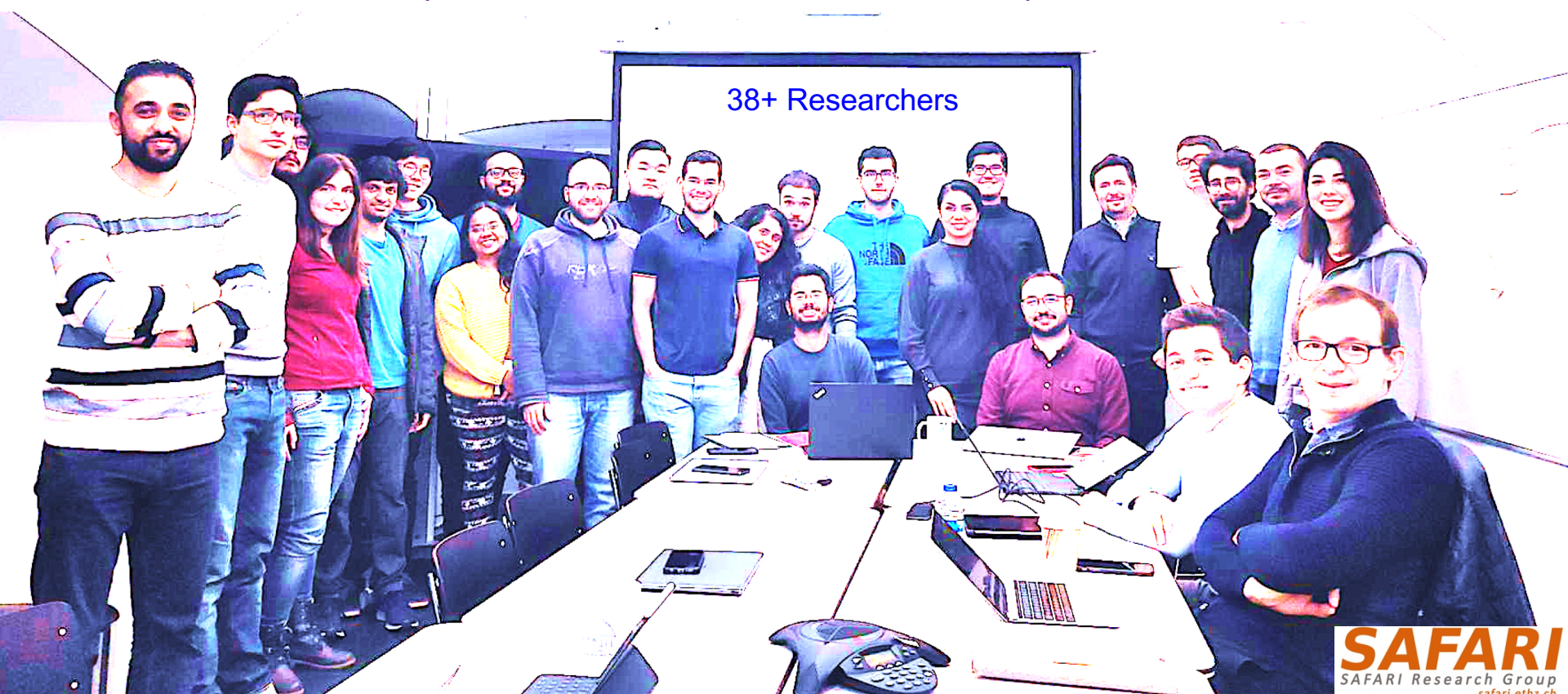
- Data-centric arch. for low energy & high perf.
  - Proc. in Mem/DRAM, NVM, unified mem/storage
- Low-latency & predictable architectures
  - Low-latency, low-energy yet low-cost memory
  - QoS-aware and predictable memory systems
- Fundamentally secure/reliable/safe arch.
  - Tolerating all bit flips; patchable HW; secure mem
- Architectures for ML/AI/Genomics/Graph/Med
  - Algorithm/arch./logic co-design; full heterogeneity
- Data-driven and data-aware architectures
  - ML/AI-driven architectural controllers and design
  - Expressive memory and expressive systems



# Onur Mutlu's SAFARI Research Group

*Computer architecture, HW/SW, systems, bioinformatics, security, memory*

<https://safari.ethz.ch/safari-newsletter-april-2020/>



# Think BIG, Aim HIGH!

**SAFARI**

<https://safari.ethz.ch>

# Principle: Teaching and Research

---

...

Teaching drives Research

Research drives Teaching

...



# Principle: Insight and Ideas

---

Focus on Insight

Encourage New Ideas

# Research & Teaching: Some Overview Talks

---

<https://www.youtube.com/onurmutlulectures>

## ■ Future Computing Architectures

- [https://www.youtube.com/watch?v=kqiZISOcGFM&list=PL5Q2soXY2Zi8D\\_5MGV6EnXEJHnV2YFBJI&index=1](https://www.youtube.com/watch?v=kqiZISOcGFM&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=1)

## ■ Enabling In-Memory Computation

- [https://www.youtube.com/watch?v=njX\\_14584Jw&list=PL5Q2soXY2Zi8D\\_5MGV6EnXEJHnV2YFBJI&index=16](https://www.youtube.com/watch?v=njX_14584Jw&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=16)

## ■ Accelerating Genome Analysis

- [https://www.youtube.com/watch?v=hPnSmfwu2-A&list=PL5Q2soXY2Zi8D\\_5MGV6EnXEJHnV2YFBJI&index=9](https://www.youtube.com/watch?v=hPnSmfwu2-A&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=9)

## ■ Rethinking Memory System Design

- [https://www.youtube.com/watch?v=F7xZLNMIY1E&list=PL5Q2soXY2Zi8D\\_5MGV6EnXEJHnV2YFBJI&index=3](https://www.youtube.com/watch?v=F7xZLNMIY1E&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=3)

## ■ Intelligent Architectures for Intelligent Machines

- [https://www.youtube.com/watch?v=n8Aj\\_A0WSq8&list=PL5Q2soXY2Zi8D\\_5MGV6EnXEJHnV2YFBJI&index=22](https://www.youtube.com/watch?v=n8Aj_A0WSq8&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=22)

# An Interview on Research and Education

---

- Computing Research and Education (@ ISCA 2019)
  - [https://www.youtube.com/watch?v=8ffSEKZhmvo&list=PL5Q2soXY2Zi\\_4oP9LdL3cc8G6NIjD2Ydz](https://www.youtube.com/watch?v=8ffSEKZhmvo&list=PL5Q2soXY2Zi_4oP9LdL3cc8G6NIjD2Ydz)
- Maurice Wilkes Award Speech (10 minutes)
  - [https://www.youtube.com/watch?v=tcQ3zZ3JpuA&list=PL5Q2soXY2Zi8D\\_5MGV6EnXEJHnV2YFBJI&index=15](https://www.youtube.com/watch?v=tcQ3zZ3JpuA&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJI&index=15)

# More Thoughts and Suggestions

---

- Onur Mutlu,  
**"Some Reflections (on DRAM)"**  
*Award Speech for ACM SIGARCH Maurice Wilkes Award, at the **ISCA** Awards Ceremony, Phoenix, AZ, USA, 25 June 2019.*  
[[Slides \(pptx\)](#) ([pdf](#))]  
[[Video of Award Acceptance Speech \(Youtube; 10 minutes\)](#) ([Youku; 13 minutes](#))]  
[[Video of Interview after Award Acceptance \(Youtube; 1 hour 6 minutes\)](#) ([Youku; 1 hour 6 minutes](#))]  
[[News Article on "ACM SIGARCH Maurice Wilkes Award goes to Prof. Onur Mutlu"](#)]
  
- Onur Mutlu,  
**"How to Build an Impactful Research Group"**  
*57th Design Automation Conference Early Career Workshop (**DAC**), Virtual, 19 July 2020.*  
[[Slides \(pptx\)](#) ([pdf](#))]

# Why Study Computer Architecture?

# Computer Architecture

---

- is the **science** and **art** of designing **computing platforms** (hardware, interface, system SW, and programming model)
- to achieve a set of **design goals**
  - E.g., highest performance on earth on workloads X, Y, Z
  - E.g., longest battery life at a form factor that fits in your pocket with cost < \$\$\$ CHF
  - E.g., best average performance across all known workloads at the best performance/cost ratio
  - ...
- Designing a supercomputer is different from designing a smartphone → But, many fundamental principles are similar

# Different Platforms, Different Goals

---



# Different Platforms, Different Goals





# Different Platforms, Different Goals

---



# Different Platforms, Different Goals

---





# Different Platforms, Different Goals





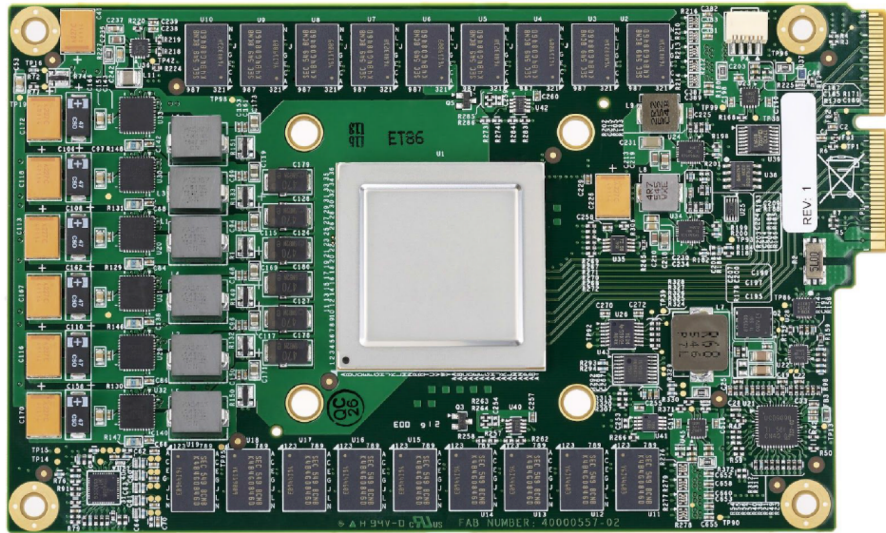
# Different Platforms, Different Goals

---

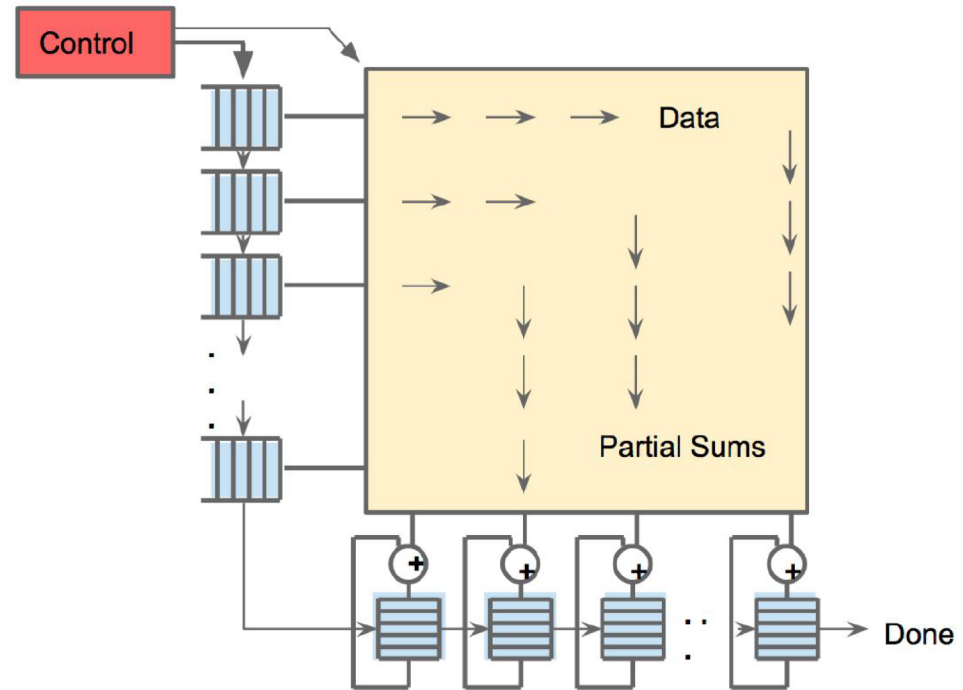


Jack Dongarra

# Different Platforms, Different Goals



**Figure 3.** TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.



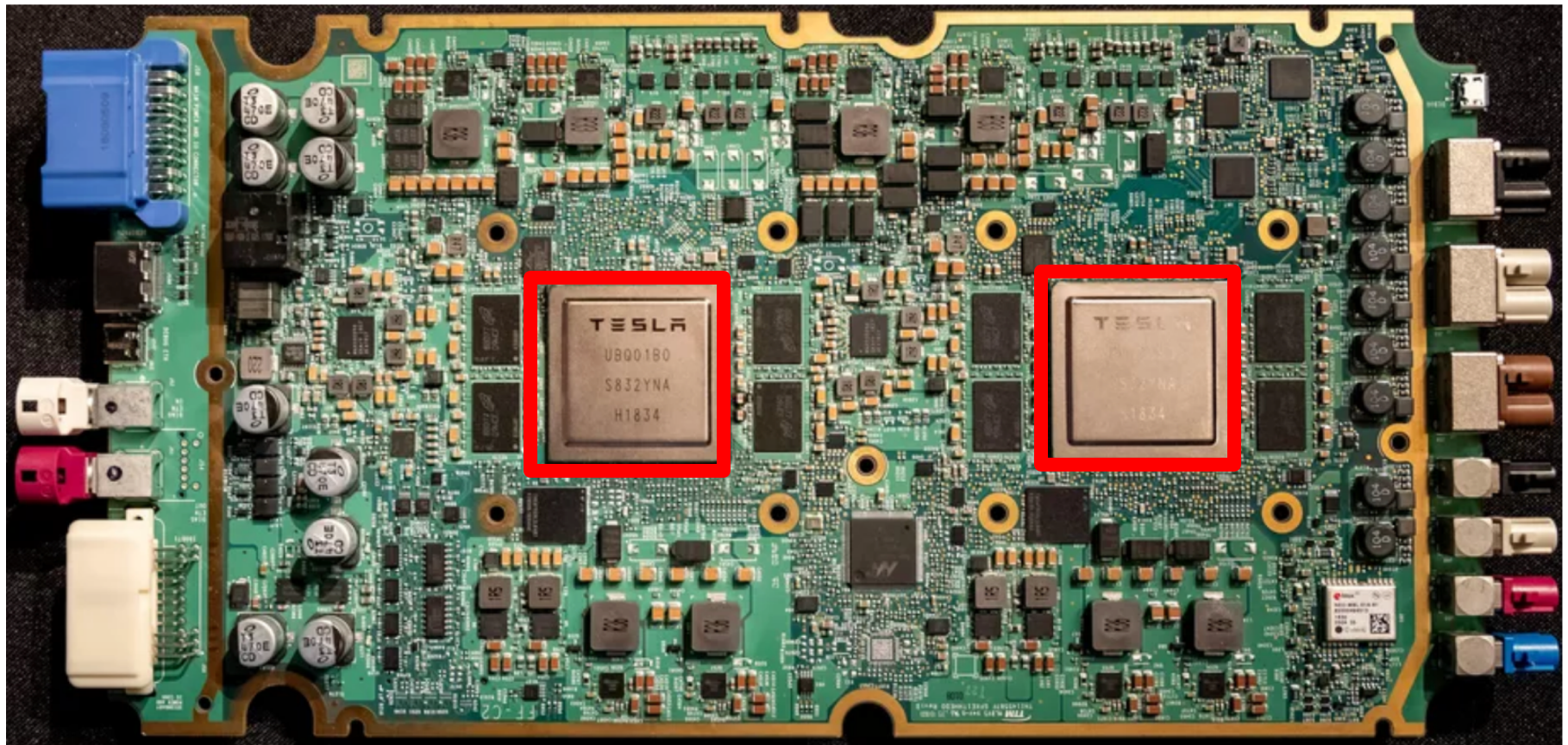
**Figure 4.** Systolic data flow of the Matrix Multiply Unit. Software has the illusion that each 256B input is read at once, and they instantly update one location of each of 256 accumulator RAMs.

Jouppi et al., “In-Datcenter Performance Analysis of a Tensor Processing Unit”, ISCA 2017.



# Different Platforms, Different Goals

- ML accelerator: 260 mm<sup>2</sup>, 6 billion transistors, 600 GFLOPS GPU, 12 ARM 2.2 GHz CPUs.
- Two redundant chips for better safety.



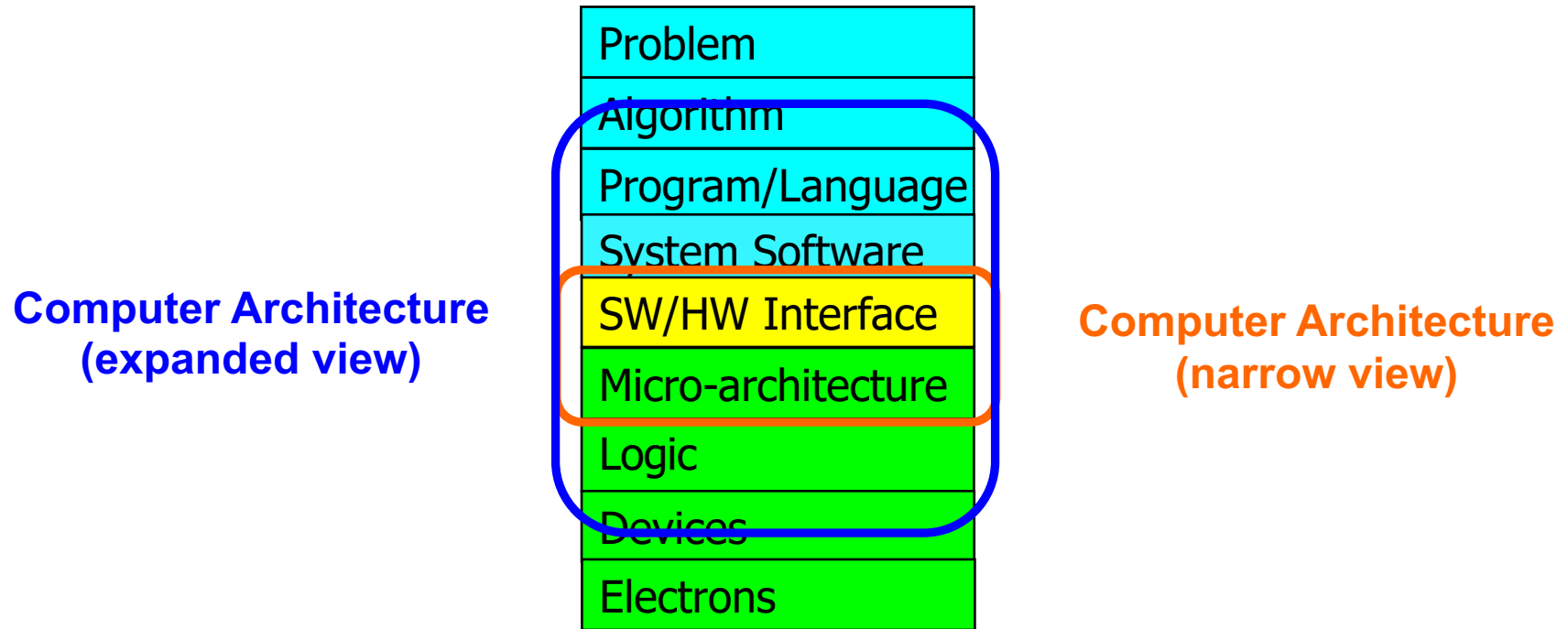
# What is Computer Architecture?

---

- The science and art of designing, selecting, and interconnecting hardware components and designing the hardware/software interface to create a computing system that meets functional, performance, energy consumption, cost, and other specific goals.

# The Transformation Hierarchy

---





# Why Study Computer Architecture?

---

- **Enable better systems:** make computers faster, cheaper, smaller, more reliable, ...
  - By exploiting advances and changes in underlying technology/circuits
- **Enable new applications**
  - Life-like 3D visualization 20 years ago? Virtual reality?
  - Self-driving cars?
  - Personalized genomics? Personalized medicine?
- **Enable better solutions** to problems
  - Software innovation is built on trends and changes in computer architecture
    - > 50% performance improvement per year has enabled this innovation
- **Understand why computers work the way they do**

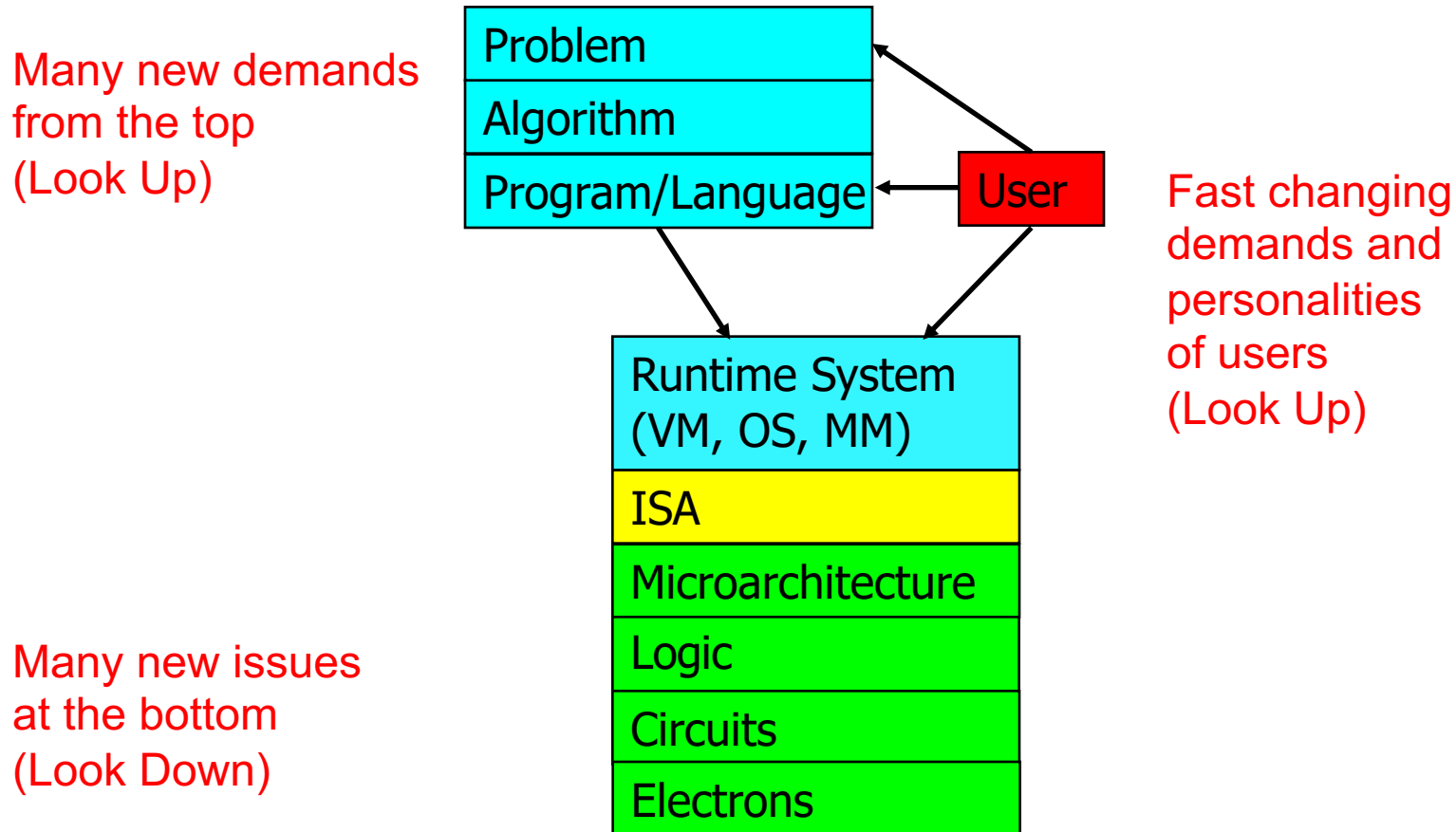
# Computer Architecture Today (I)

---

- Today is a very exciting time to study computer architecture
  - Industry is in a large paradigm shift (to novel architectures)
    - many different potential system designs possible
  - **Many difficult problems** *motivating* and *caused by* the shift
    - Huge hunger for data and new data-intensive applications
    - Power/energy/thermal constraints
    - Complexity of design
    - Difficulties in technology scaling
    - Memory bottleneck
    - Reliability problems
    - Programmability problems
    - Security and privacy issues
  - No clear, definitive answers to these problems
-

# Computer Architecture Today (II)

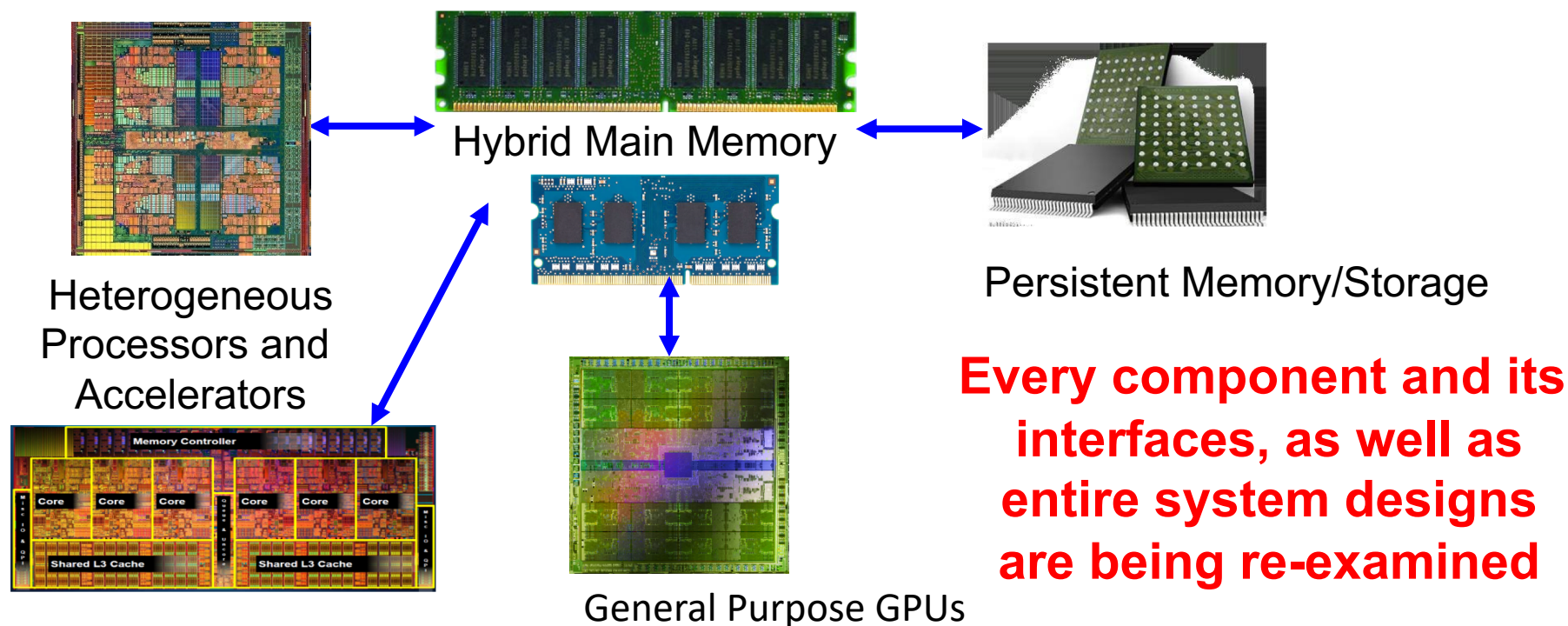
- These problems affect all parts of the computing stack – if we do not change the way we design systems



- No clear, definitive answers to these problems

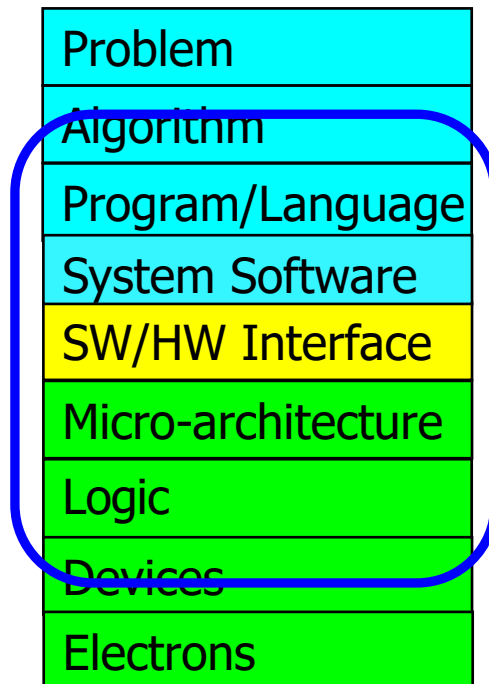
# Computer Architecture Today (III)

- Computing landscape is very different from 10-20 years ago
- Both UP (software and humanity trends) and DOWN (technologies and their issues), FORWARD and BACKWARD, and the resulting requirements and constraints



To achieve the highest **energy efficiency** and **performance**:

**we must take the expanded view**  
of computer architecture



**Co-design across the hierarchy:**  
**Algorithms to devices**

**Specialize as much as possible**  
**within the design goals**

# Historical: Opportunities at the Bottom

---

## There's Plenty of Room at the Bottom

---

From Wikipedia, the free encyclopedia

**"There's Plenty of Room at the Bottom: An Invitation to Enter a New Field of Physics"** was a lecture given by [physicist Richard Feynman](#) at the annual [American Physical Society](#) meeting at [Caltech](#) on December 29, 1959.<sup>[1]</sup> Feynman considered the possibility of direct manipulation of individual atoms as a more powerful form of synthetic chemistry than those used at the time. Although versions of the talk were reprinted in a few popular magazines, it went largely unnoticed and did not inspire the conceptual beginnings of the field. Beginning in the 1980s, nanotechnology advocates cited it to establish the scientific credibility of their work.

# Historical: Opportunities at the Bottom (II)

---

## There's Plenty of Room at the Bottom

---

From Wikipedia, the free encyclopedia

Feynman considered some ramifications of a general ability to manipulate matter on an atomic scale. He was particularly interested in the possibilities of denser [computer](#) circuitry, and [microscopes](#) that could see things much smaller than is possible with [scanning electron microscopes](#). These ideas were later realized by the use of the [scanning tunneling microscope](#), the [atomic force microscope](#) and other examples of [scanning probe microscopy](#) and storage systems such as [Millipede](#), created by researchers at [IBM](#).

Feynman also suggested that it should be possible, in principle, to make [nanoscale machines](#) that "arrange the atoms the way we want", and do chemical synthesis by mechanical manipulation.

He also presented the possibility of "[swallowing the doctor](#)", an idea that he credited in the essay to his friend and graduate student [Albert Hibbs](#). This concept involved building a tiny, swallowable surgical robot.



# Historical: Opportunities at the Top

## REVIEW

### There's plenty of room at the Top: What will drive computer performance after Moore's law?

 Charles E. Leiserson<sup>1</sup>,  Neil C. Thompson<sup>1,2,\*</sup>,  Joel S. Emer<sup>1,3</sup>,  Bradley C. Kuszmaul<sup>1,†</sup>, Butler W. Lampson<sup>1,4</sup>,  ...

+ See all authors and affiliations

*Science* 05 Jun 2020:  
Vol. 368, Issue 6495, eaam9744  
DOI: 10.1126/science.aam9744

Much of the improvement in computer performance comes from decades of miniaturization of computer components, a trend that was foreseen by the Nobel Prize-winning physicist Richard Feynman in his 1959 address, “There’s Plenty of Room at the Bottom,” to the American Physical Society. In 1975, Intel founder Gordon Moore predicted the regularity of this miniaturization trend, now called Moore’s law, which, until recently, doubled the number of transistors on computer chips every 2 years.

Unfortunately, semiconductor miniaturization is running out of steam as a viable way to grow computer performance—there isn’t much more room at the “Bottom.” If growth in computing power stalls, practically all industries will face challenges to their productivity. Nevertheless, opportunities for growth in computing performance will still be available, especially at the “Top” of the computing-technology stack: software, algorithms, and hardware architecture.



# Axiom, Revisited

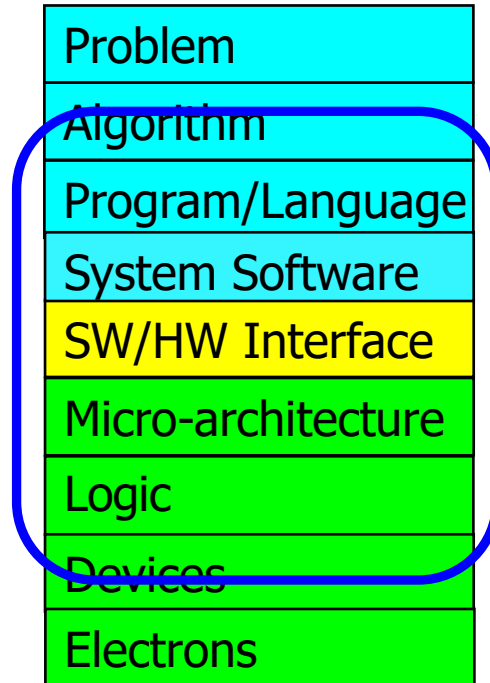
---

There is plenty of room both at the top and at the bottom  
but much more so  
when you communicate well between and optimize across  
the top and the bottom.

# Hence the Expanded View

---

**Computer Architecture  
(expanded view)**



# Some Cross-Layer Design Examples (Foreshadowing)

# Expressive (Memory) Interfaces

---

- Nandita Vijaykumar, Abhilasha Jain, Diptesh Majumdar, Kevin Hsieh, Gennady Pekhimenko, Eiman Ebrahimi, Nastaran Hajinazar, Phillip B. Gibbons and Onur Mutlu, **"A Case for Richer Cross-layer Abstractions: Bridging the Semantic Gap with Expressive Memory"**  
*Proceedings of the 45th International Symposium on Computer Architecture (ISCA)*, Los Angeles, CA, USA, June 2018.  
[[Slides \(pptx\)](#)] [[pdf](#)] [[Lightning Talk Slides \(pptx\)](#)] [[pdf](#)]  
[[Lightning Talk Video](#)]

---

## A Case for Richer Cross-layer Abstractions: Bridging the Semantic Gap with Expressive Memory

Nandita Vijaykumar<sup>†§</sup> Abhilasha Jain<sup>†</sup> Diptesh Majumdar<sup>†</sup> Kevin Hsieh<sup>†</sup> Gennady Pekhimenko<sup>‡</sup>  
Eiman Ebrahimi<sup>⌘</sup> Nastaran Hajinazar<sup>†</sup> Phillip B. Gibbons<sup>†</sup> Onur Mutlu<sup>§†</sup>

<sup>†</sup>Carnegie Mellon University

<sup>‡</sup>University of Toronto

<sup>⌘</sup>NVIDIA

<sup>†</sup>Simon Fraser University

<sup>§</sup>ETH Zürich

# X-MeM Aids Many Optimizations

**Table 1: Summary of the example memory optimizations that XMem aids.**

Memory optimization	Example semantics provided by XMem (described in §3.3)	Example Benefits of XMem
Cache management	(i) Distinguishing between data structures or pools of similar data; (ii) Working set size; (iii) Data reuse	Enables: (i) applying different caching policies to different data structures or pools of data; (ii) avoiding cache thrashing by <i>knowing</i> the active working set size; (iii) bypassing/prioritizing data that has no/high reuse. (§5)
Page placement in DRAM e.g., [23, 24]	(i) Distinguishing between data structures; (ii) Access pattern; (iii) Access intensity	Enables page placement at the <i>data structure</i> granularity to (i) isolate data structures that have high row buffer locality and (ii) spread out concurrently-accessed irregular data structures across banks and channels to improve parallelism. (§6)
Cache/memory compression e.g., [25–32]	(i) Data type: integer, float, char; (ii) Data properties: sparse, pointer, data index	Enables using a <i>different compression algorithm</i> for each data structure based on data type and data properties, e.g., sparse data encodings, FP-specific compression, delta-based compression for pointers [27].
Data prefetching e.g., [33–36]	(i) Access pattern: strided, irregular, irregular but repeated (e.g., graphs), access stride; (ii) Data type: index, pointer	Enables (i) <i>highly accurate</i> software-driven prefetching while leveraging the benefits of hardware prefetching (e.g., by being memory bandwidth-aware, avoiding cache thrashing); (ii) using different prefetcher <i>types</i> for different data structures: e.g., stride [33], tile-based [20], pattern-based [34–37], data-based for indices/pointers [38, 39], etc.
DRAM cache management e.g., [40–46]	(i) Access intensity; (ii) Data reuse; (iii) Working set size	(i) Helps avoid cache thrashing by knowing working set size [44]; (ii) Better DRAM cache management via reuse behavior and access intensity information.
Approximation in memory e.g., [47–53]	(i) Distinguishing between pools of similar data; (ii) Data properties: tolerance towards approximation	Enables (i) each memory component to track how approximable data is (at a fine granularity) to inform approximation techniques; (ii) data placement in heterogeneous reliability memories [54].
Data placement: NUMA systems e.g., [55, 56]	(i) Data partitioning across threads (i.e., relating data to threads that access it); (ii) Read-Write properties	Reduces the need for profiling or data migration (i) to co-locate data with threads that access it and (ii) to identify Read-Only data, thereby enabling techniques such as replication.
Data placement: hybrid memories e.g., [16, 57, 58]	(i) Read-Write properties (Read-Only/Read-Write); (ii) Access intensity; (iii) Data structure size; (iv) Access pattern	Avoids the need for profiling/migration of data in hybrid memories to (i) effectively manage the asymmetric read-write properties in NVM (e.g., placing Read-Only data in the NVM) [16, 57]; (ii) make tradeoffs between data structure "hotness" and size to allocate fast/high bandwidth memory [14]; and (iii) leverage row-buffer locality in placement based on access pattern [45].
Managing NUCA systems e.g., [15, 59]	(i) Distinguishing pools of similar data; (ii) Access intensity; (iii) Read-Write or Private-Shared properties	(i) Enables using different cache policies for different data pools (similar to [15]); (ii) Reduces the need for reactive mechanisms that detect sharing and read-write characteristics to inform cache policies.



# Expressive (Memory) Interfaces for GPUs

---

- Nandita Vijaykumar, Eiman Ebrahimi, Kevin Hsieh, Phillip B. Gibbons and Onur Mutlu, **"The Locality Descriptor: A Holistic Cross-Layer Abstraction to Express Data Locality in GPUs"**  
*Proceedings of the 45th International Symposium on Computer Architecture (ISCA)*, Los Angeles, CA, USA, June 2018.  
[[Slides \(pptx\)](#)] [[pdf](#)] [[Lightning Talk Slides \(pptx\)](#)] [[pdf](#)]  
[[Lightning Talk Video](#)]

## The Locality Descriptor: A Holistic Cross-Layer Abstraction to Express Data Locality in GPUs

Nandita Vijaykumar <sup>†§</sup>	Eiman Ebrahimi <sup>‡</sup>	Kevin Hsieh <sup>†</sup>
Phillip B. Gibbons <sup>†</sup>	Onur Mutlu <sup>§†</sup>	
<sup>†</sup> Carnegie Mellon University	<sup>‡</sup> NVIDIA	<sup>§</sup> ETH Zürich

# Heterogeneous-Reliability Memory

---

- Yixin Luo, Sriram Govindan, Bikash Sharma, Mark Santaniello, Justin Meza, Aman Kansal, Jie Liu, Badriddine Khessib, Kushagra Vaid, and Onur Mutlu,  
**"Characterizing Application Memory Error Vulnerability to Optimize Data Center Cost via Heterogeneous-Reliability Memory"**  
*Proceedings of the 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, Atlanta, GA, June 2014. [[Summary](#)]  
[[Slides \(pptx\)](#)] [[pdf](#)] [[Coverage on ZDNet](#)]

## Characterizing Application Memory Error Vulnerability to Optimize Datacenter Cost via Heterogeneous-Reliability Memory

Yixin Luo   Sriram Govindan\*   Bikash Sharma\*   Mark Santaniello\*   Justin Meza  
Aman Kansal\*   Jie Liu\*   Badriddine Khessib\*   Kushagra Vaid\*   Onur Mutlu

Carnegie Mellon University, yixinluo@cs.cmu.edu, {meza, onur}@cmu.edu

\*Microsoft Corporation, {srgovin, bsharma, marksan, kansal, jie.liu, bknessib, kvaid}@microsoft.com

# EDEN: Data-Aware Efficient DNN Inference

---

- Skanda Koppula, Lois Orosa, A. Giray Yaglikci, Roknoddin Azizi, Taha Shahroodi, Konstantinos Kanellopoulos, and Onur Mutlu,  
**"EDEN: Enabling Energy-Efficient, High-Performance Deep Neural Network Inference Using Approximate DRAM"**  
*Proceedings of the 52nd International Symposium on Microarchitecture (MICRO)*, Columbus, OH, USA, October 2019.  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Lightning Talk Slides \(pptx\)](#)] [[pdf](#)]  
[[Poster \(pptx\)](#)] [[pdf](#)]  
[[Lightning Talk Video](#) (90 seconds)]  
[[Full Talk Lecture](#) (38 minutes)]

## EDEN: Enabling Energy-Efficient, High-Performance Deep Neural Network Inference Using Approximate DRAM

Skanda Koppula   Lois Orosa   A. Giray Yağlıkçı  
Roknoddin Azizi   Taha Shahroodi   Konstantinos Kanellopoulos   Onur Mutlu

ETH Zürich

# SMASH: SW/HW Indexing Acceleration

---

- Konstantinos Kanellopoulos, Nandita Vijaykumar, Christina Giannoula, Roknoddin Azizi, Skanda Koppula, Nika Mansouri Ghiasi, Taha Shahroodi, Juan Gomez-Luna, and Onur Mutlu,

## **"SMASH: Co-designing Software Compression and Hardware-Accelerated Indexing for Efficient Sparse Matrix Operations"**

*Proceedings of the 52nd International Symposium on Microarchitecture (MICRO), Columbus, OH, USA, October 2019.*

[[Slides \(pptx\)](#) ([pdf](#))]

[[Lightning Talk Slides \(pptx\)](#) ([pdf](#))]

[[Poster \(pptx\)](#) ([pdf](#))]

[[Lightning Talk Video](#) (90 seconds)]

[[Full Talk Lecture](#) (30 minutes)]

## **SMASH: Co-designing Software Compression and Hardware-Accelerated Indexing for Efficient Sparse Matrix Operations**

Konstantinos Kanellopoulos<sup>1</sup> Nandita Vijaykumar<sup>2,1</sup> Christina Giannoula<sup>1,3</sup> Roknoddin Azizi<sup>1</sup>  
Skanda Koppula<sup>1</sup> Nika Mansouri Ghiasi<sup>1</sup> Taha Shahroodi<sup>1</sup> Juan Gomez Luna<sup>1</sup> Onur Mutlu<sup>1,2</sup>

<sup>1</sup>ETH Zürich

<sup>2</sup>Carnegie Mellon University

<sup>3</sup>National Technical University of Athens

# Rethinking Virtual Memory

---

- Nastaran Hajinazar, Pratyush Patel, Minesh Patel, Konstantinos Kanellopoulos, Saugata Ghose, Rachata Ausavarungnirun, Geraldo Francisco de Oliveira Jr., Jonathan Appavoo, Vivek Seshadri, and Onur Mutlu,  
**"The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework"**  
*Proceedings of the 47th International Symposium on Computer Architecture (ISCA)*, Valencia, Spain, June 2020.  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Lightning Talk Slides \(pptx\)](#)] [[pdf](#)]  
[[ARM Research Summit Poster \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#) (26 minutes)]  
[[Lightning Talk Video](#) (3 minutes)]

## The Virtual Block Interface: A Flexible Alternative to the Conventional Virtual Memory Framework

Nastaran Hajinazar<sup>\*†</sup> Pratyush Patel<sup>✕</sup> Minesh Patel<sup>\*</sup> Konstantinos Kanellopoulos<sup>\*</sup> Saugata Ghose<sup>‡</sup>  
Rachata Ausavarungnirun<sup>⊙</sup> Geraldo F. Oliveira<sup>\*</sup> Jonathan Appavoo<sup>◇</sup> Vivek Seshadri<sup>▽</sup> Onur Mutlu<sup>\*‡</sup>

<sup>\*</sup>ETH Zürich   <sup>†</sup>Simon Fraser University   <sup>✕</sup>University of Washington   <sup>‡</sup>Carnegie Mellon University  
<sup>⊙</sup>King Mongkut's University of Technology North Bangkok   <sup>◇</sup>Boston University   <sup>▽</sup>Microsoft Research India



# Many Interesting Things Are Happening Today in Computer Architecture

Many Interesting Things  
Are Happening Today  
in Computer Architecture

**Performance  
and  
Energy Efficiency**

# Intel Optane Persistent Memory (2019)

---

- Non-volatile main memory
- Based on 3D-XPoint Technology



# PCM as Main Memory: Idea in 2009

---

- Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger, **"Architecting Phase Change Memory as a Scalable DRAM Alternative"**  
*Proceedings of the 36th International Symposium on Computer Architecture (ISCA)*, pages 2-13, Austin, TX, June 2009. [Slides \(pdf\)](#)

## Architecting Phase Change Memory as a Scalable DRAM Alternative

Benjamin C. Lee<sup>†</sup> Engin Ipek<sup>†</sup> Onur Mutlu<sup>‡</sup> Doug Burger<sup>†</sup>

<sup>†</sup>Computer Architecture Group  
Microsoft Research  
Redmond, WA  
{blee, ipek, dburger}@microsoft.com

<sup>‡</sup>Computer Architecture Laboratory  
Carnegie Mellon University  
Pittsburgh, PA  
onur@cmu.edu

# PCM as Main Memory: Idea in 2009

---

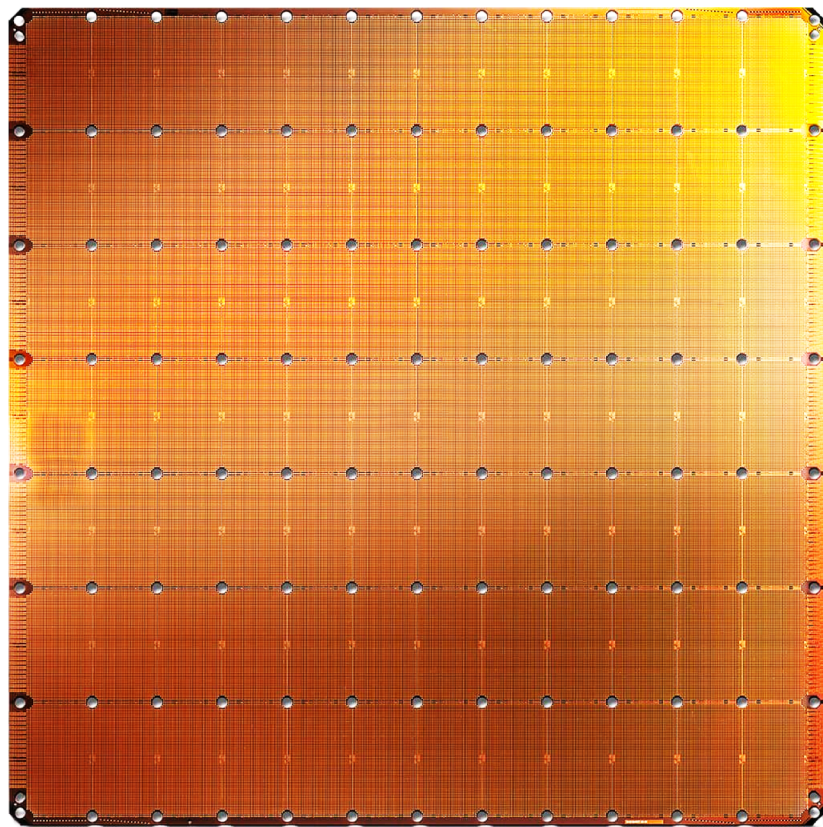
- Benjamin C. Lee, Ping Zhou, Jun Yang, Youtao Zhang, Bo Zhao, Engin Ipek, Onur Mutlu, and Doug Burger,  
**"Phase Change Technology and the Future of Main Memory"**  
*IEEE Micro, Special Issue: Micro's Top Picks from 2009 Computer Architecture Conferences (**MICRO TOP PICKS**), Vol. 30, No. 1, pages 60-70, January/February 2010.*

## PHASE-CHANGE TECHNOLOGY AND THE FUTURE OF MAIN MEMORY



# Cerebras's Wafer Scale Engine (2019)

---



## **Cerebras WSE**

1.2 Trillion transistors

46,225 mm<sup>2</sup>

- The largest ML accelerator chip
- 400,000 cores



## **Largest GPU**

21.1 Billion transistors

815 mm<sup>2</sup>

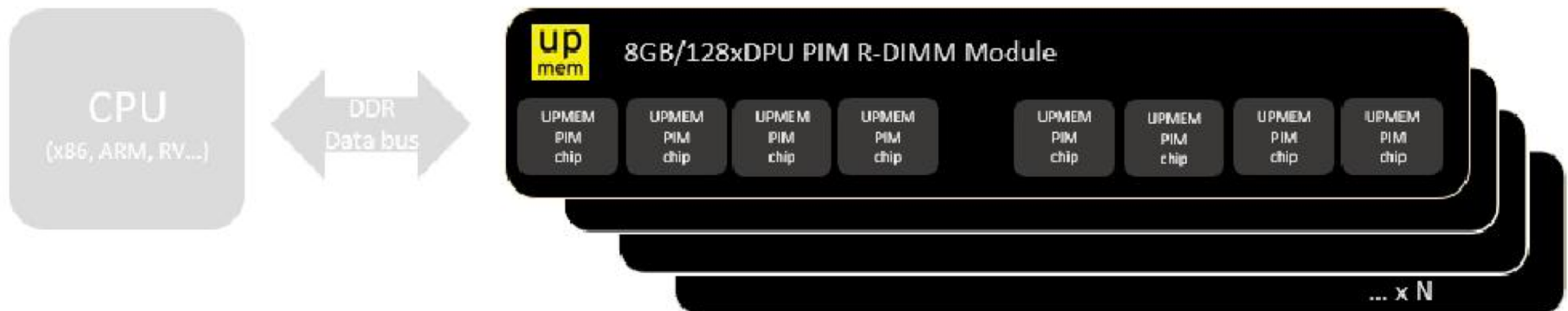
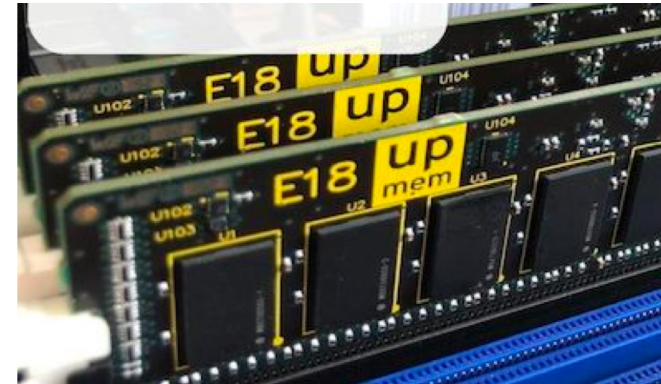
NVIDIA TITAN V

<https://www.anandtech.com/show/14758/hot-chips-31-live-blogs-cerebras-wafer-scale-deep-learning>

<https://www.cerebras.net/cerebras-wafer-scale-engine-why-we-need-big-chips-for-deep-learning/>

# UPMEM Processing-in-DRAM Engine (2019)

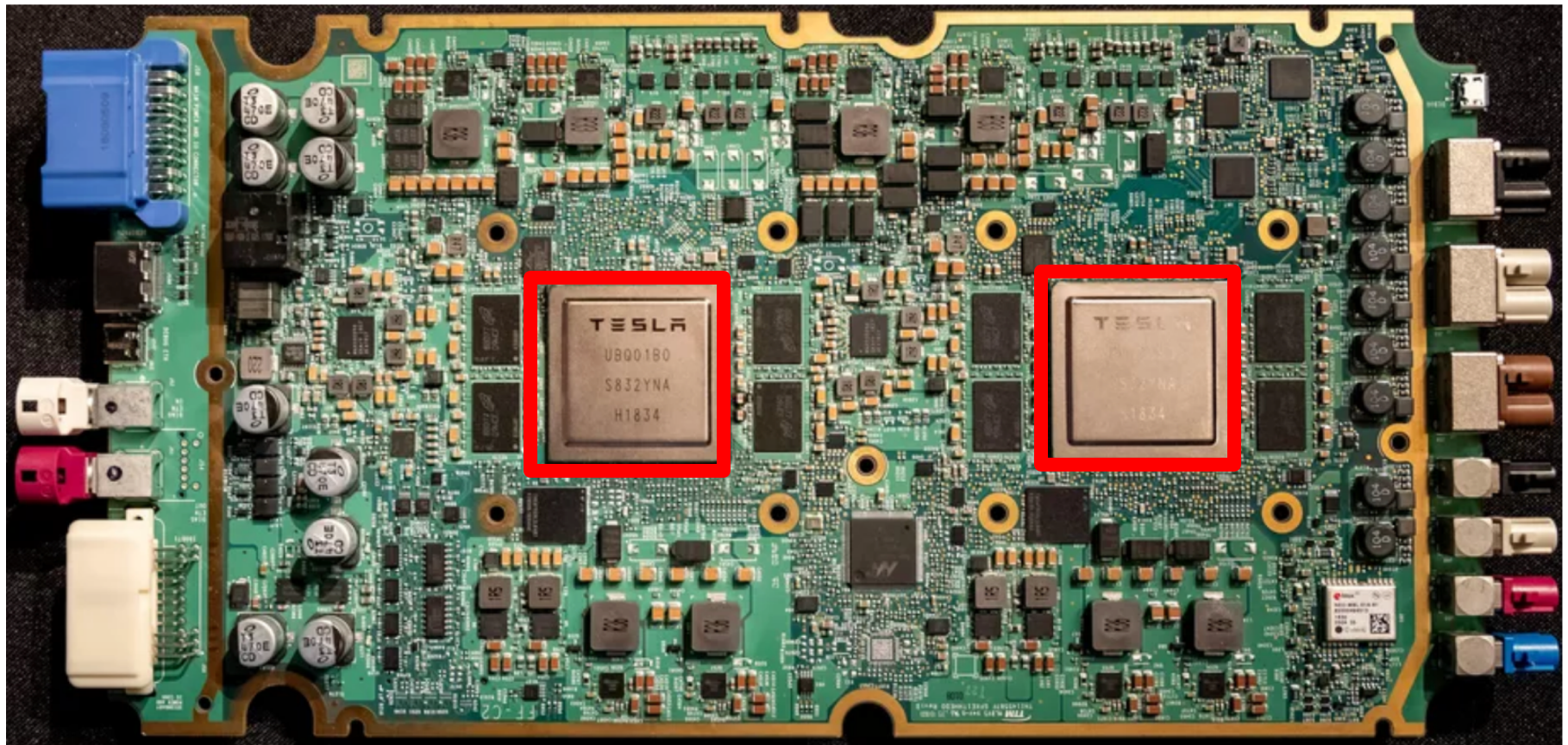
- **Processing in DRAM Engine**
- Includes **standard DIMM modules**, with a **large number of DPU processors** combined with DRAM chips.
- Replaces **standard DIMMs**
  - DDR4 R-DIMM modules
    - 8GB+128 DPUs (16 PIM chips)
    - Standard 2x-nm DRAM process
  - **Large amounts of** compute & memory bandwidth



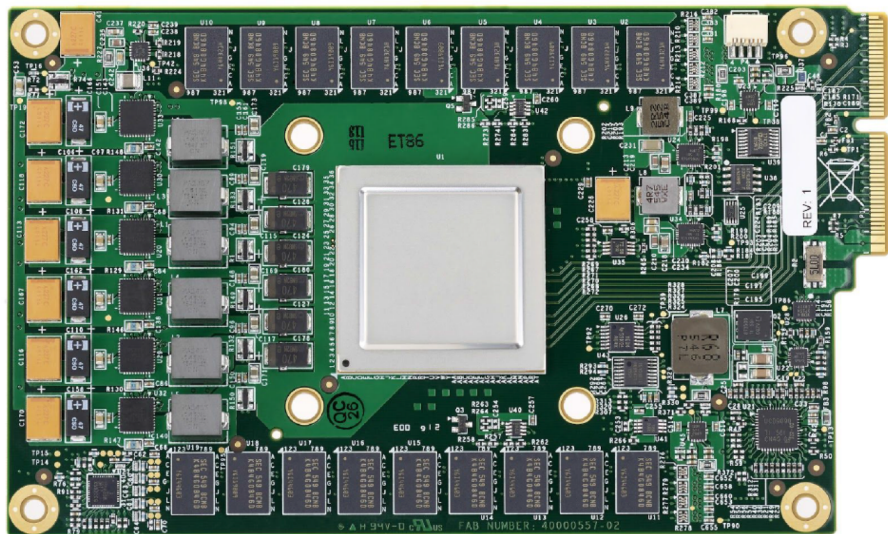


# TESLA Full Self-Driving Computer (2019)

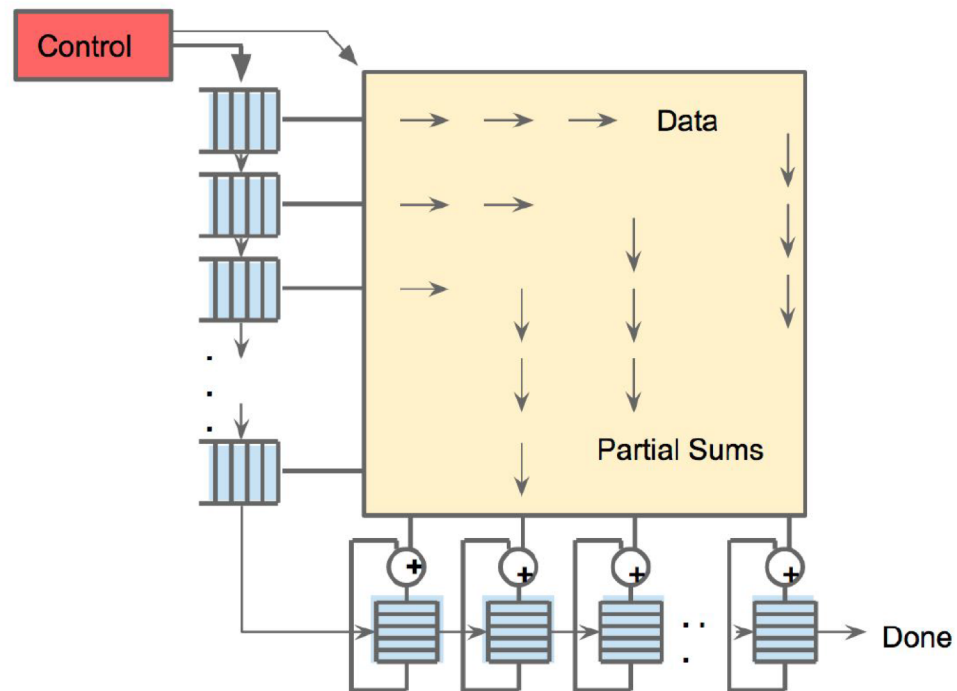
- ML accelerator: 260 mm<sup>2</sup>, 6 billion transistors, 600 GFLOPS GPU, 12 ARM 2.2 GHz CPUs.
- Two redundant chips for better safety.



# Google TPU Generation I (~2016)



**Figure 3.** TPU Printed Circuit Board. It can be inserted in the slot for an SATA disk in a server, but the card uses PCIe Gen3 x16.



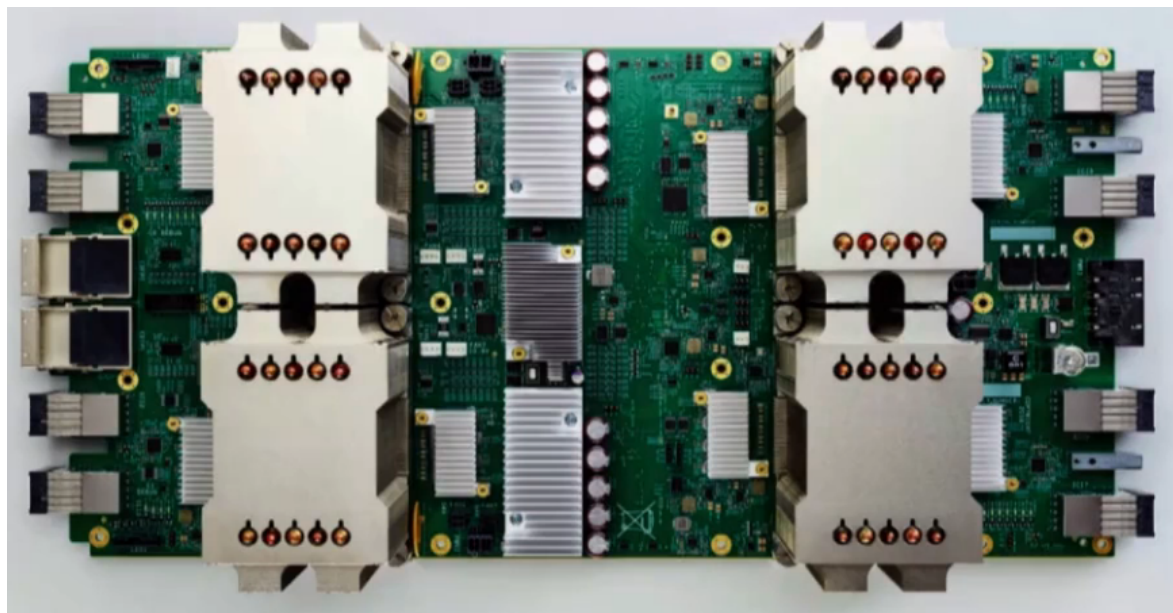
**Figure 4.** Systolic data flow of the Matrix Multiply Unit. Software has the illusion that each 256B input is read at once, and they instantly update one location of each of 256 accumulator RAMs.

Jouppi et al., “In-Datacenter Performance Analysis of a Tensor Processing Unit”, ISCA 2017.



# Google TPU Generation II (2017)

---



<https://www.nextplatform.com/2017/05/17/first-depth-look-googles-new-second-generation-tpu/>

4 TPU chips  
vs 1 chip in TPU1

High Bandwidth Memory  
vs DDR3

Floating point operations  
vs FP16

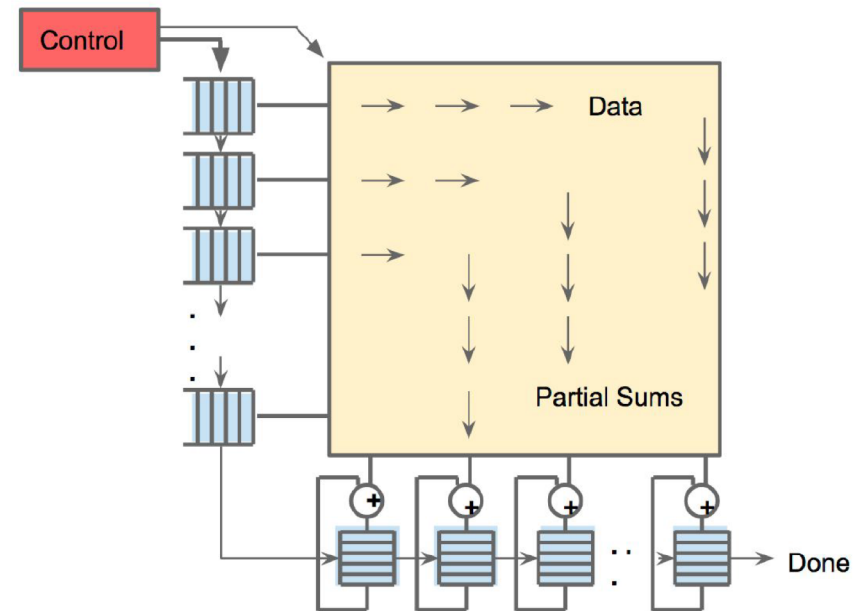
45 TFLOPS per chip  
vs 23 TOPS

Designed for training  
and inference  
vs only inference



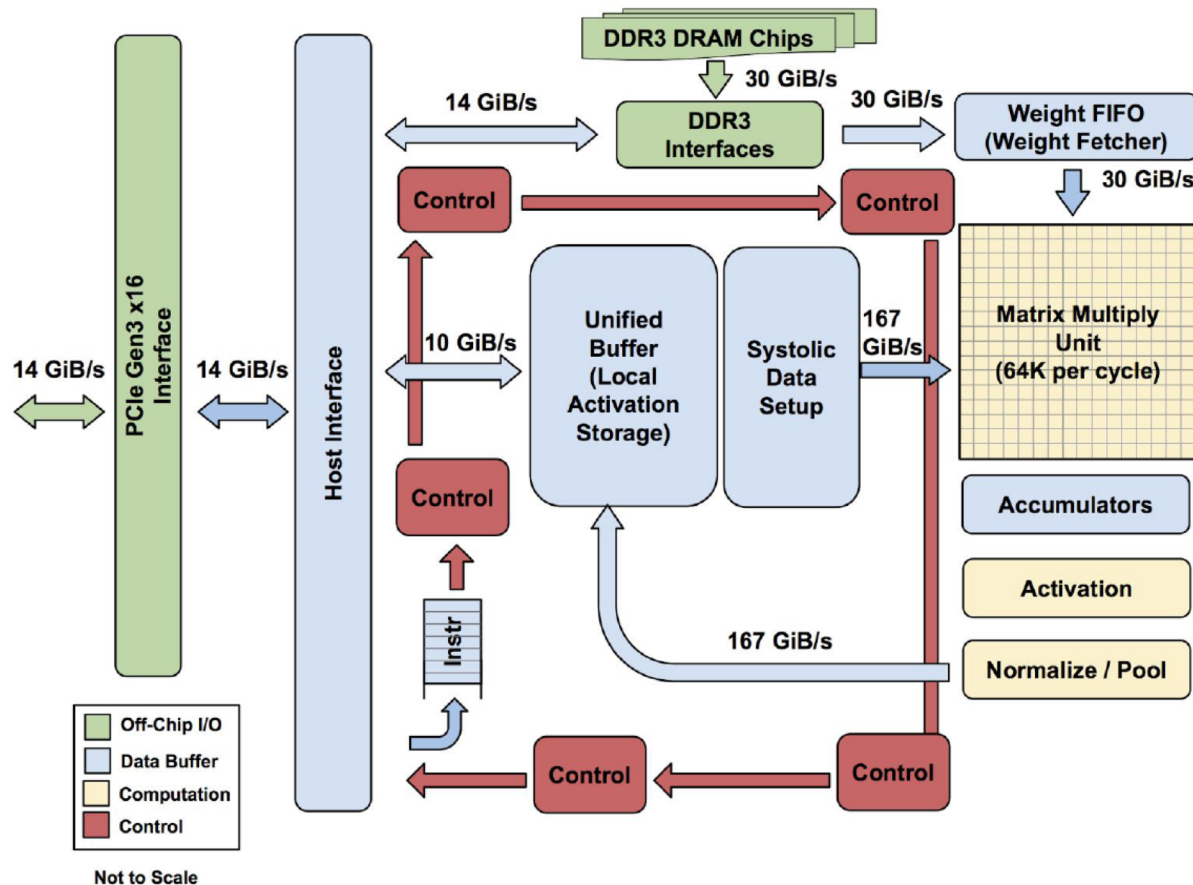
# An Example Modern Systolic Array: TPU (II)

As reading a large SRAM uses much more power than arithmetic, the matrix unit uses systolic execution to save energy by reducing reads and writes of the Unified Buffer [Kun80][Ram91][Ovt15b]. Figure 4 shows that data flows in from the left, and the weights are loaded from the top. A given 256-element multiply-accumulate operation moves through the matrix as a diagonal wavefront. The weights are preloaded, and take effect with the advancing wave alongside the first data of a new block. Control and data are pipelined to give the illusion that the 256 inputs are read at once, and that they instantly update one location of each of 256 accumulators. From a correctness perspective, software is unaware of the systolic nature of the matrix unit, but for performance, it does worry about the latency of the unit.



Jouppi et al., “In-Datacenter Performance Analysis of a Tensor Processing Unit”, ISCA 2017.

# An Example Modern Systolic Array: TPU (III)



**Figure 1.** TPU Block Diagram. The main computation part is the yellow Matrix Multiply unit in the upper right hand corner. Its inputs are the blue Weight FIFO and the blue Unified Buffer (UB) and its output is the blue Accumulators (Acc). The yellow Activation Unit performs the nonlinear functions on the Acc, which go to the UB.

# Many (Other) AI/ML Chips

---

- Alibaba
- Amazon
- Facebook
- Google
- Huawei
- Intel
- Microsoft
- NVIDIA
- Tesla
- Many Others and Many Startups...
- **Many More to Come...**

# Many (Other) AI/ML Chips

## AI Chip Landscape

S.T.



All information contained within this infographic is gathered from the internet and periodically updated, no guarantee is given that the information provided is correct, complete, and up-to-date.

# Many Interesting Things Are Happening Today in Computer Architecture

Many Interesting Things  
Are Happening Today  
in Computer Architecture

**Reliability  
and  
Security**



# Security: RowHammer (2014)



# The Story of RowHammer

- One can **predictably induce bit flips** in commodity DRAM chips
  - >80% of the tested DRAM chips are vulnerable
- First example of how a **simple hardware failure mechanism** can create a **widespread system security vulnerability**

**WIRED**

Forget Software—Now Hackers Are Exploiting Physics

BUSINESS	CULTURE	DESIGN	GEAR	SCIENCE
----------	---------	--------	------	---------

ANDY GREENBERG SECURITY 08.31.16 7:00 AM

SHARE



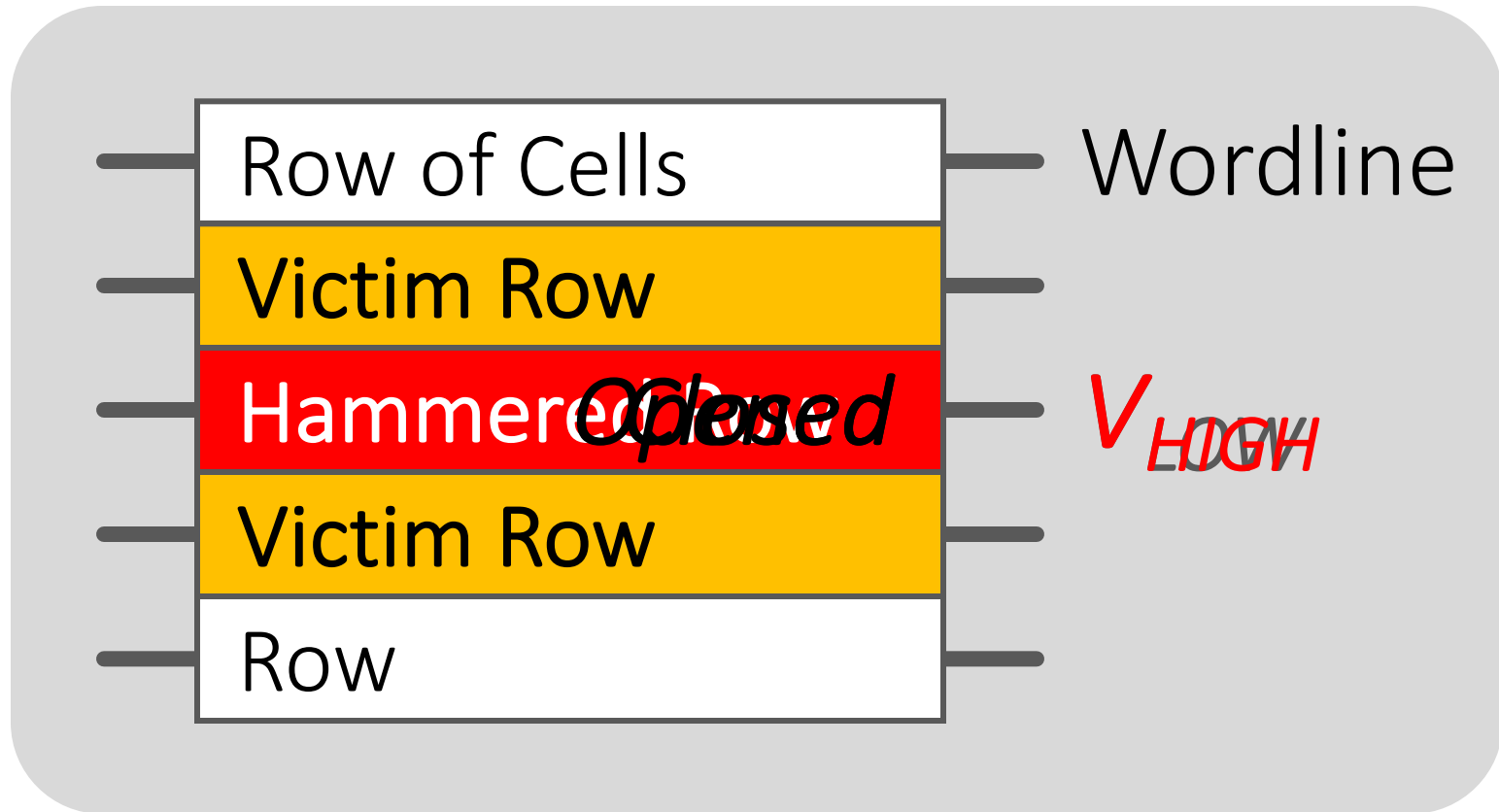
SHARE  
18276



TWEET

# FORGET SOFTWARE—NOW HACKERS ARE EXPLOITING PHYSICS

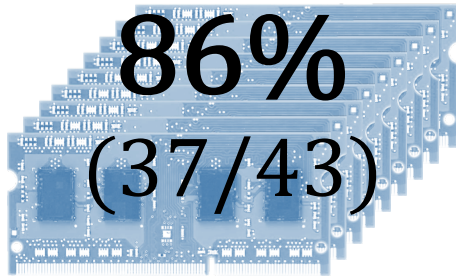
# Modern DRAM is Prone to Disturbance Errors



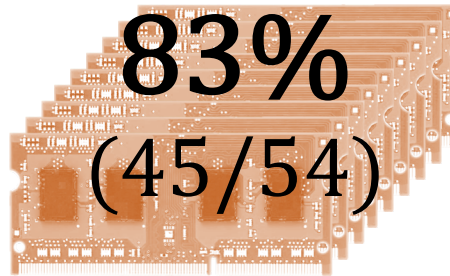
Repeatedly reading a row enough times (before memory gets refreshed) induces **disturbance errors** in adjacent rows in **most real DRAM chips you can buy today**

# Most DRAM Modules Are Vulnerable

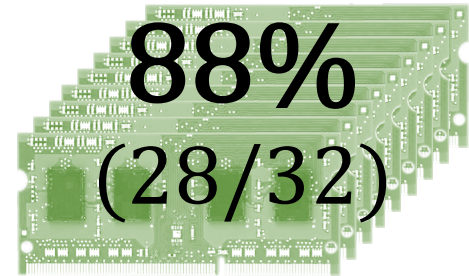
A company



B company



C company



Up to  
 $1.0 \times 10^7$   
errors

Up to  
 $2.7 \times 10^6$   
errors

Up to  
 $3.3 \times 10^5$   
errors

# One Can Take Over an Otherwise-Secure System

---

## Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors

*Abstract. Memory isolation is a key property of a reliable and secure computing system — an access to one memory address should not have unintended side effects on data stored in other addresses. However, as DRAM process technology*

## Project Zero

Flipping Bits in Memory Without Accessing Them:  
An Experimental Study of DRAM Disturbance Errors  
(Kim et al., ISCA 2014)

News and updates from the Project Zero team at Google

Exploiting the DRAM rowhammer bug to  
gain kernel privileges (Seaborn, 2015)

Monday, March 9, 2015

Exploiting the DRAM rowhammer bug to gain kernel privileges



# Security: RowHammer (2014)



It's like breaking into an apartment by repeatedly slamming a neighbor's door until the vibrations open the door you were after



# RowHammer: Five Years Ago...

---

- Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu,  
**"Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors"**  
*Proceedings of the 41st International Symposium on Computer Architecture (ISCA)*, Minneapolis, MN, June 2014.  
[[Slides \(pptx\)](#)] [[pdf](#)] [[Lightning Session Slides \(pptx\)](#)] [[pdf](#)] [[Source Code and Data](#)]

## Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors

Yoongu Kim<sup>1</sup>   Ross Daly\*   Jeremie Kim<sup>1</sup>   Chris Fallin\*   Ji Hye Lee<sup>1</sup>  
Donghyuk Lee<sup>1</sup>   Chris Wilkerson<sup>2</sup>   Konrad Lai   Onur Mutlu<sup>1</sup>

<sup>1</sup>Carnegie Mellon University   <sup>2</sup>Intel Labs

# RowHammer: Now and Beyond...

---

- Onur Mutlu and Jeremie Kim,  
**"RowHammer: A Retrospective"**  
*IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD) Special Issue on Top Picks in Hardware and Embedded Security*, 2019.  
[[Preliminary arXiv version](#)]

## RowHammer: A Retrospective

Onur Mutlu<sup>§‡</sup>      Jeremie S. Kim<sup>‡§</sup>  
§ETH Zürich      ‡Carnegie Mellon University

# RowHammer in 2020 (I)

---

- Jeremie S. Kim, Minesh Patel, A. Giray Yaglikci, Hasan Hassan, Roknoddin Azizi, Lois Orosa, and Onur Mutlu,  
**"Revisiting RowHammer: An Experimental Analysis of Modern Devices and Mitigation Techniques"**  
*Proceedings of the 47th International Symposium on Computer Architecture (ISCA)*, Valencia, Spain, June 2020.  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Lightning Talk Slides \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#) (20 minutes)]  
[[Lightning Talk Video](#) (3 minutes)]

## Revisiting RowHammer: An Experimental Analysis of Modern DRAM Devices and Mitigation Techniques

Jeremie S. Kim<sup>§†</sup>      Minesh Patel<sup>§</sup>      A. Giray Yağlıkçı<sup>§</sup>  
Hasan Hassan<sup>§</sup>      Roknoddin Azizi<sup>§</sup>      Lois Orosa<sup>§</sup>      Onur Mutlu<sup>§†</sup>  
<sup>§</sup>*ETH Zürich*      <sup>†</sup>*Carnegie Mellon University*

# RowHammer in 2020 (II)

---

- Pietro Frigo, Emanuele Vannacci, Hasan Hassan, Victor van der Veen, Onur Mutlu, Cristiano Giuffrida, Herbert Bos, and Kaveh Razavi, **"TRRespass: Exploiting the Many Sides of Target Row Refresh"**  
*Proceedings of the 41st IEEE Symposium on Security and Privacy (S&P)*, San Francisco, CA, USA, May 2020.  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#) (17 minutes)]  
[[Source Code](#)]  
[[Web Article](#)]  
***Best paper award.***

## TRRespass: Exploiting the Many Sides of Target Row Refresh

Pietro Frigo<sup>\*†</sup>   Emanuele Vannacci<sup>\*†</sup>   Hasan Hassan<sup>§</sup>   Victor van der Veen<sup>¶</sup>  
Onur Mutlu<sup>§</sup>   Cristiano Giuffrida<sup>\*</sup>   Herbert Bos<sup>\*</sup>   Kaveh Razavi<sup>\*</sup>

# RowHammer in 2020 (III)

---

- Lucian Cojocar, Jeremie Kim, Minesh Patel, Lillian Tsai, Stefan Saroiu, Alec Wolman, and Onur Mutlu,  
**"Are We Susceptible to Rowhammer? An End-to-End Methodology for Cloud Providers"**  
*Proceedings of the 41st IEEE Symposium on Security and Privacy (S&P)*, San Francisco, CA, USA, May 2020.  
[[Slides \(pptx\)](#)] [[pdf](#)]  
[[Talk Video](#) (17 minutes)]

## Are We Susceptible to Rowhammer?

## An End-to-End Methodology for Cloud Providers

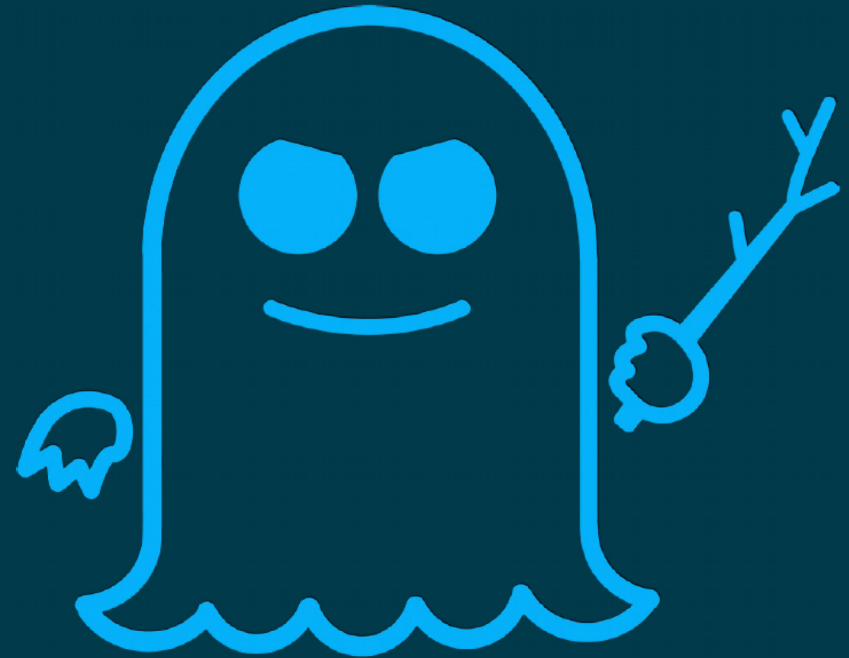
Lucian Cojocar, Jeremie Kim<sup>§†</sup>, Minesh Patel<sup>§</sup>, Lillian Tsai<sup>‡</sup>,  
Stefan Saroiu, Alec Wolman, and Onur Mutlu<sup>§†</sup>  
Microsoft Research, <sup>§</sup>ETH Zürich, <sup>†</sup>CMU, <sup>‡</sup>MIT

# Security: Meltdown and Spectre (2018)

---



**MELTDOWN**



**SPECTRE**



# Meltdown and Spectre

---

- Someone can steal secret data from the system even though
  - ❑ your program and data are perfectly correct and
  - ❑ your hardware behaves according to the specification and
  - ❑ there are no software vulnerabilities/bugs
  
- Why?
  - ❑ Speculative execution leaves traces of secret data in the processor's cache (internal storage)
    - It brings data that is not supposed to be brought/accessed if there was no speculative execution
  - ❑ A malicious program can inspect the contents of the cache to "infer" secret data that it is not supposed to access
  - ❑ A malicious program can actually force another program to speculatively execute code that leaves traces of secret data

# More on Meltdown/Spectre Vulnerabilities

---

## Project Zero

News and updates from the Project Zero team at Google

Wednesday, January 3, 2018

### Reading privileged memory with a side-channel

Posted by Jann Horn, Project Zero

We have discovered that CPU data cache timing can be abused to efficiently leak information out of mis-speculated execution, leading to (at worst) arbitrary virtual memory read vulnerabilities across local security boundaries in various contexts.

# Many Interesting Things Are Happening Today in Computer Architecture

Many Interesting Things  
Are Happening Today  
in Computer Architecture

**More Demanding Workloads**

# New Genome Sequencing Technologies

---

## Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions

Damla Senol Cali ✉, Jeremie S Kim, Saugata Ghose, Can Alkan, Onur Mutlu

*Briefings in Bioinformatics*, bby017, <https://doi.org/10.1093/bib/bby017>

**Published:** 02 April 2018    **Article history** ▼



Oxford Nanopore MinION

Data → performance & energy bottleneck

# Why Do We Care? An Example

200 Oxford Nanopore sequencers have left UK for China, to support rapid, near-sample coronavirus sequencing for outbreak surveillance

Fri 31st January 2020

Following extensive support of, and collaboration with, public health professionals in China, Oxford Nanopore has shipped an additional 200 MinION sequencers and related consumables to China. These will be used to support the ongoing surveillance of the current coronavirus outbreak, adding to a large number of the devices already installed in the country.

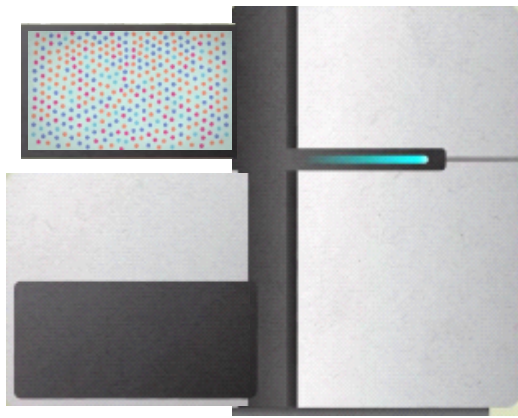


Each MinION sequencer is approximately the size of a stapler, and can provide rapid sequence information about the coronavirus.



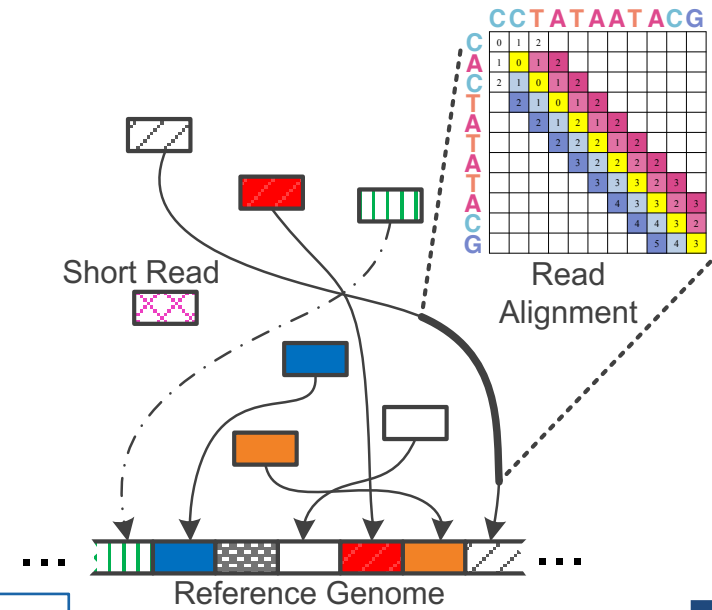
700Kg of Oxford Nanopore sequencers and consumables are on their way for use by Chinese scientists in understanding the current coronavirus outbreak.





Billions of Short Reads

ATATATACGTACTAGTACGT  
 TTTAGTACGTACGT  
 ATACGTACTAGTACGT  
 CGCCCCTACGTA  
 ACGTACTAGTACGT  
 TTAGTACGTACGT  
 TACGTACTAAAGTACGT  
 TACGTACTAGTACGT  
 TTTAAACGTA  
 CGTACTAGTACGT  
 GGGAGTACGTACGT



## 1 Sequencing

# Genome Analysis

## 2 Read Mapping

Data → performance & energy bottleneck

read4: CGCTTCCAT  
 read5: CCATGACGC  
 read6: TTCCATGAC



## 3 Variant Calling

## 4 Scientific Discovery

# GateKeeper: FPGA-Based Alignment Filtering

---

- Mohammed Alser, Hasan Hassan, Hongyi Xin, Oguz Ergin, Onur Mutlu, and Can Alkan  
**"GateKeeper: A New Hardware Architecture for Accelerating Pre-Alignment in DNA Short Read Mapping"**  
**Bioinformatics**, [published online, May 31], 2017.  
[[Source Code](#)]  
[[Online link at Bioinformatics Journal](#)]

## GateKeeper: a new hardware architecture for accelerating pre-alignment in DNA short read mapping

Mohammed Alser ✉, Hasan Hassan, Hongyi Xin, Oğuz Ergin, Onur Mutlu ✉, Can Alkan ✉

*Bioinformatics*, Volume 33, Issue 21, 1 November 2017, Pages 3355–3363,

<https://doi.org/10.1093/bioinformatics/btx342>

**Published:** 31 May 2017    **Article history** ▼

# In-Memory DNA Sequence Analysis

---

- Jeremie S. Kim, Damla Senol Cali, Hongyi Xin, Donghyuk Lee, Saugata Ghose, Mohammed Alser, Hasan Hassan, Oguz Ergin, Can Alkan, and Onur Mutlu, **"GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies"** ***BMC Genomics***, 2018.  
*Proceedings of the 16th Asia Pacific Bioinformatics Conference (APBC)*, Yokohama, Japan, January 2018.  
[arxiv.org Version \(pdf\)](https://arxiv.org/abs/1801.00000)

## GRIM-Filter: Fast seed location filtering in DNA read mapping using processing-in-memory technologies

Jeremie S. Kim<sup>1,6\*</sup>, Damla Senol Cali<sup>1</sup>, Hongyi Xin<sup>2</sup>, Donghyuk Lee<sup>3</sup>, Saugata Ghose<sup>1</sup>, Mohammed Alser<sup>4</sup>, Hasan Hassan<sup>6</sup>, Oguz Ergin<sup>5</sup>, Can Alkan<sup>4\*</sup> and Onur Mutlu<sup>6,1\*</sup>

From The Sixteenth Asia Pacific Bioinformatics Conference 2018  
Yokohama, Japan. 15-17 January 2018

# GenASM: Fast Approximate String Matching

---

## GenASM: A High-Performance, Low-Power Approximate String Matching Acceleration Framework for Genome Sequence Analysis

Damla Senol Cali<sup>†⌘</sup> Gurpreet S. Kalsi<sup>⌘</sup> Zülal Bingöl<sup>∇</sup> Can Firtina<sup>◇</sup> Lavanya Subramanian<sup>‡</sup> Jeremie S. Kim<sup>◇†</sup>  
Rachata Ausavarungnirun<sup>⊙</sup> Mohammed Alser<sup>◇</sup> Juan Gomez-Luna<sup>◇</sup> Amirali Boroumand<sup>†</sup> Anant Nori<sup>⌘</sup>  
Allison Scibisz<sup>†</sup> Sreenivas Subramoney<sup>⌘</sup> Can Alkan<sup>∇</sup> Saugata Ghose<sup>\*†</sup> Onur Mutlu<sup>◇†∇</sup>  
<sup>†</sup>*Carnegie Mellon University*   <sup>⌘</sup>*Processor Architecture Research Lab, Intel Labs*   <sup>∇</sup>*Bilkent University*   <sup>◇</sup>*ETH Zürich*  
<sup>‡</sup>*Facebook*   <sup>⊙</sup>*King Mongkut's University of Technology North Bangkok*   <sup>\*</sup>*University of Illinois at Urbana–Champaign*

# More on Genome Analysis: Another Lecture

---

- Onur Mutlu,  
**"Accelerating Genome Analysis: A Primer on an Ongoing Journey"**  
*Keynote talk at 2nd Workshop on Accelerator Architecture in Computational Biology and Bioinformatics (AACBB), Washington, DC, USA, February 2019.*  
[[Slides \(pptx\)\(pdf\)](#)]  
[[Video](#)]

## Accelerating Genome Analysis A Primer on an Ongoing Journey

Onur Mutlu

[omutlu@gmail.com](mailto:omutlu@gmail.com)

<https://people.inf.ethz.ch/omutlu>

16 February 2019

AACBB Keynote Talk

**SAFARI**

**ETH** zürich

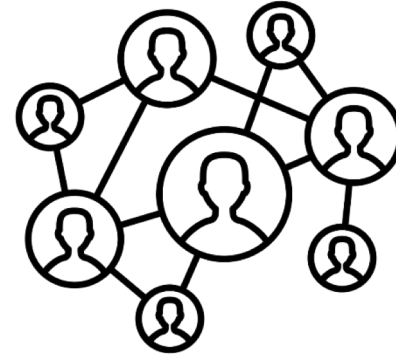
**Carnegie Mellon**

# Data Overwhelms Modern Machines

---



**In-memory Databases**



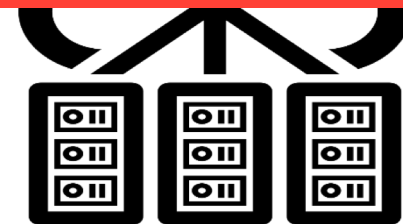
**Graph/Tree Processing**

**Data → performance & energy bottleneck**



**In-Memory Data Analytics**

[Clapp+ (Intel), IISWC'15;  
Awan+, BDCloud'15]



**Datacenter Workloads**

[Kanev+ (Google), ISCA'15]



# Data Overwhelms Modern Machines



**Chrome**



**TensorFlow Mobile**

Data → performance & energy bottleneck

**VP9**



**Video Playback**

Google's **video codec**

**VP9**



**Video Capture**

Google's **video codec**

# Data Movement Overwhelms Modern Machines

- Amirali Boroumand, Saugata Ghose, Youngsok Kim, Rachata Ausavarungnirun, Eric Shiu, Rahul Thakur, Daehyun Kim, Aki Kuusela, Allan Knies, Parthasarathy Ranganathan, and Onur Mutlu, **"Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks"** *Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, Williamsburg, VA, USA, March 2018.

**62.7% of the total system energy  
is spent on data movement**

## Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks

Amirali Boroumand<sup>1</sup>

Saugata Ghose<sup>1</sup>

Youngsok Kim<sup>2</sup>

Rachata Ausavarungnirun<sup>1</sup>

Eric Shiu<sup>3</sup>

Rahul Thakur<sup>3</sup>

Daehyun Kim<sup>4,3</sup>

Aki Kuusela<sup>3</sup>

Allan Knies<sup>3</sup>

Parthasarathy Ranganathan<sup>3</sup>

Onur Mutlu<sup>5,1</sup>

# Many Interesting Things Are Happening Today in Computer Architecture

# Many Novel Concepts Investigated Today

---

- **New Computing Paradigms (Rethinking the Full Stack)**
  - ❑ Processing in Memory, Processing Near Data
  - ❑ Neuromorphic Computing
  - ❑ Fundamentally Secure and Dependable Computers
- **New Accelerators (Algorithm-Hardware Co-Designs)**
  - ❑ Artificial Intelligence & Machine Learning
  - ❑ Graph Analytics
  - ❑ Genome Analysis
- **New Memories and Storage Systems**
  - ❑ Non-Volatile Main Memory
  - ❑ Intelligent Memory

# Increasingly Demanding Applications

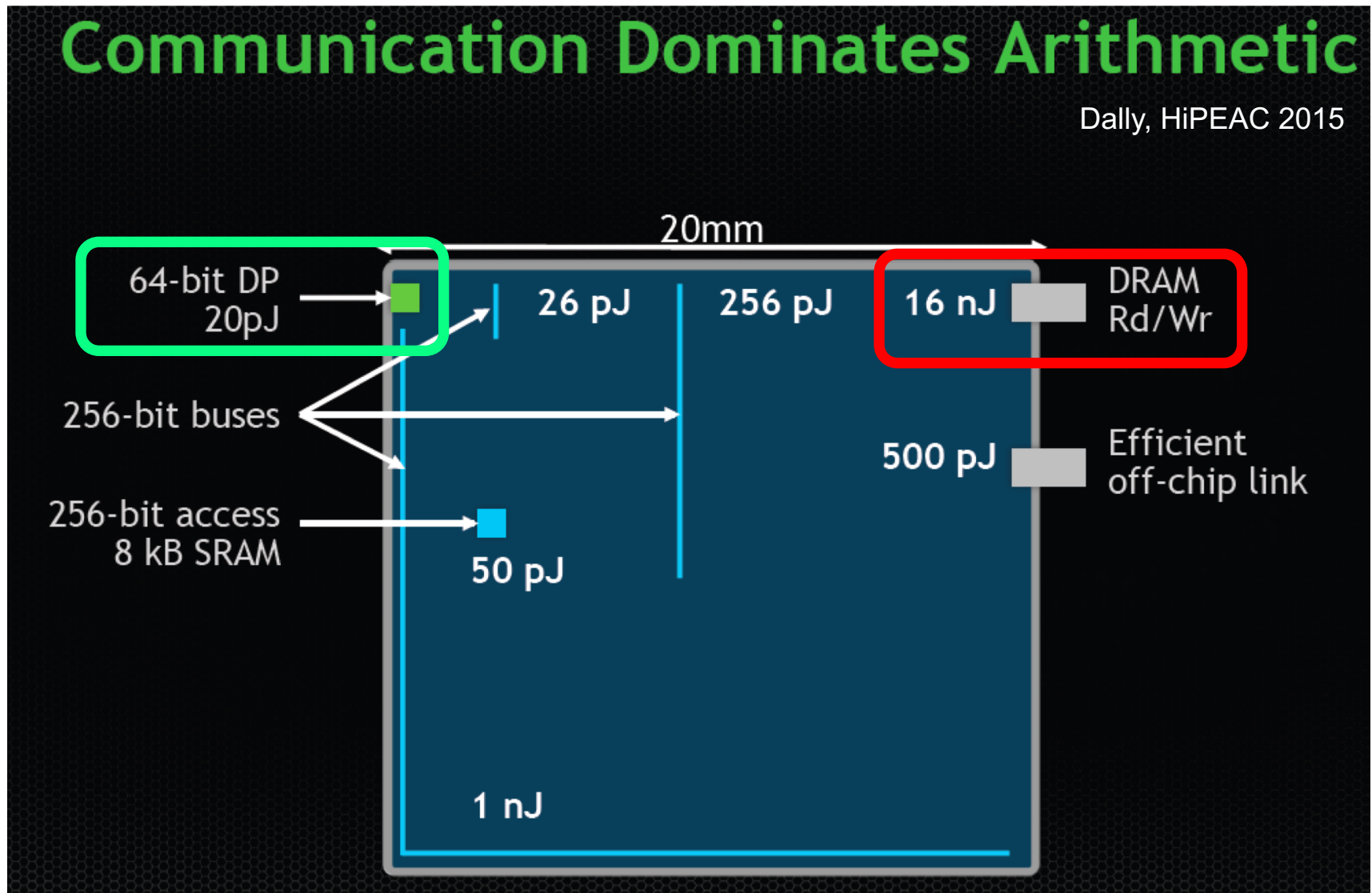
---

- Dream, and they will come

# Increasingly Diverging/Complex Tradeoffs

## Communication Dominates Arithmetic

Dally, HiPEAC 2015

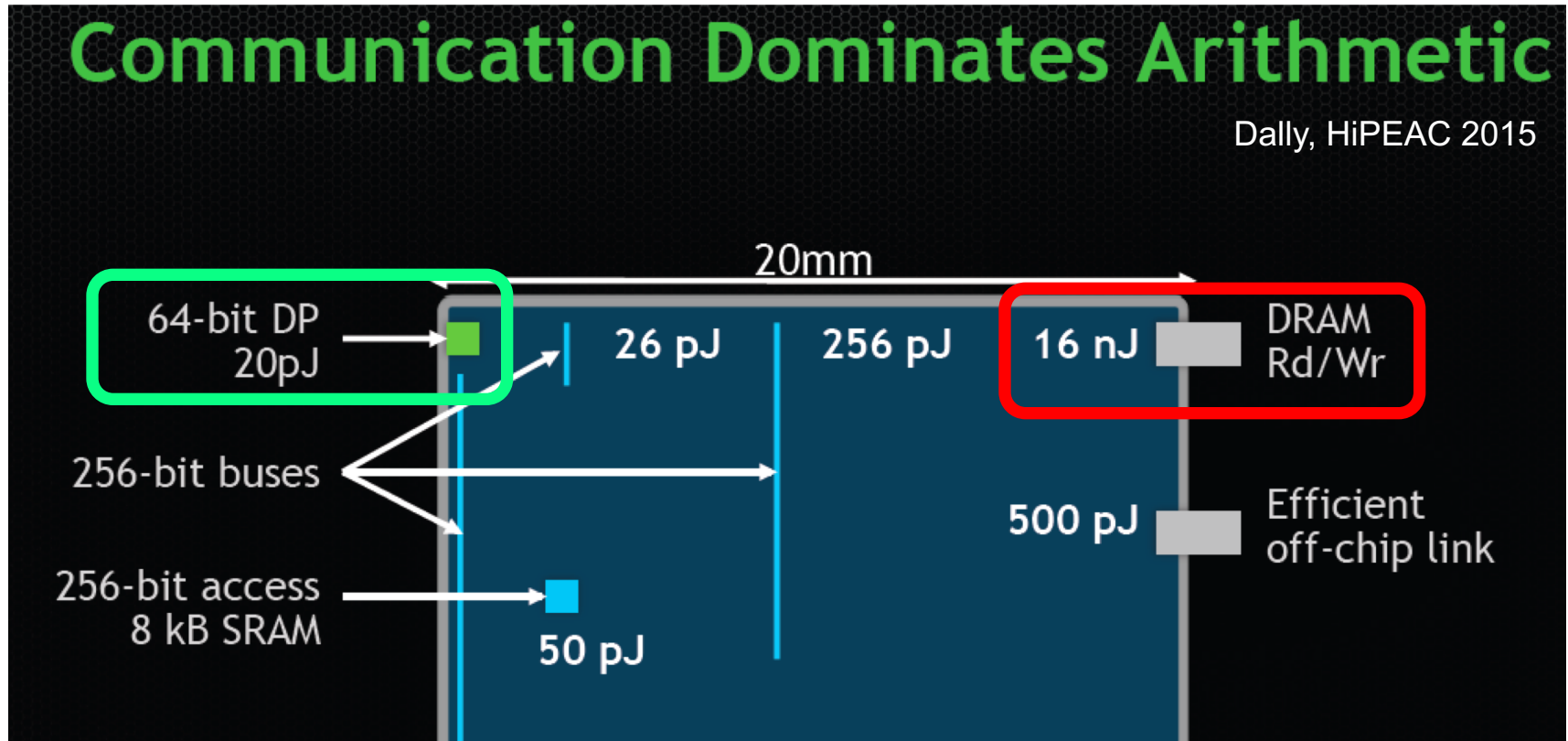




# Increasingly Diverging/Complex Tradeoffs

## Communication Dominates Arithmetic

Dally, HiPEAC 2015

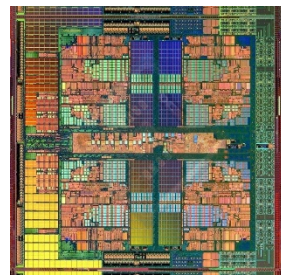


A memory access consumes  $\sim 1000\times$  the energy of a complex addition

# Increasingly Complex Systems

---

## Past systems



Microprocessor



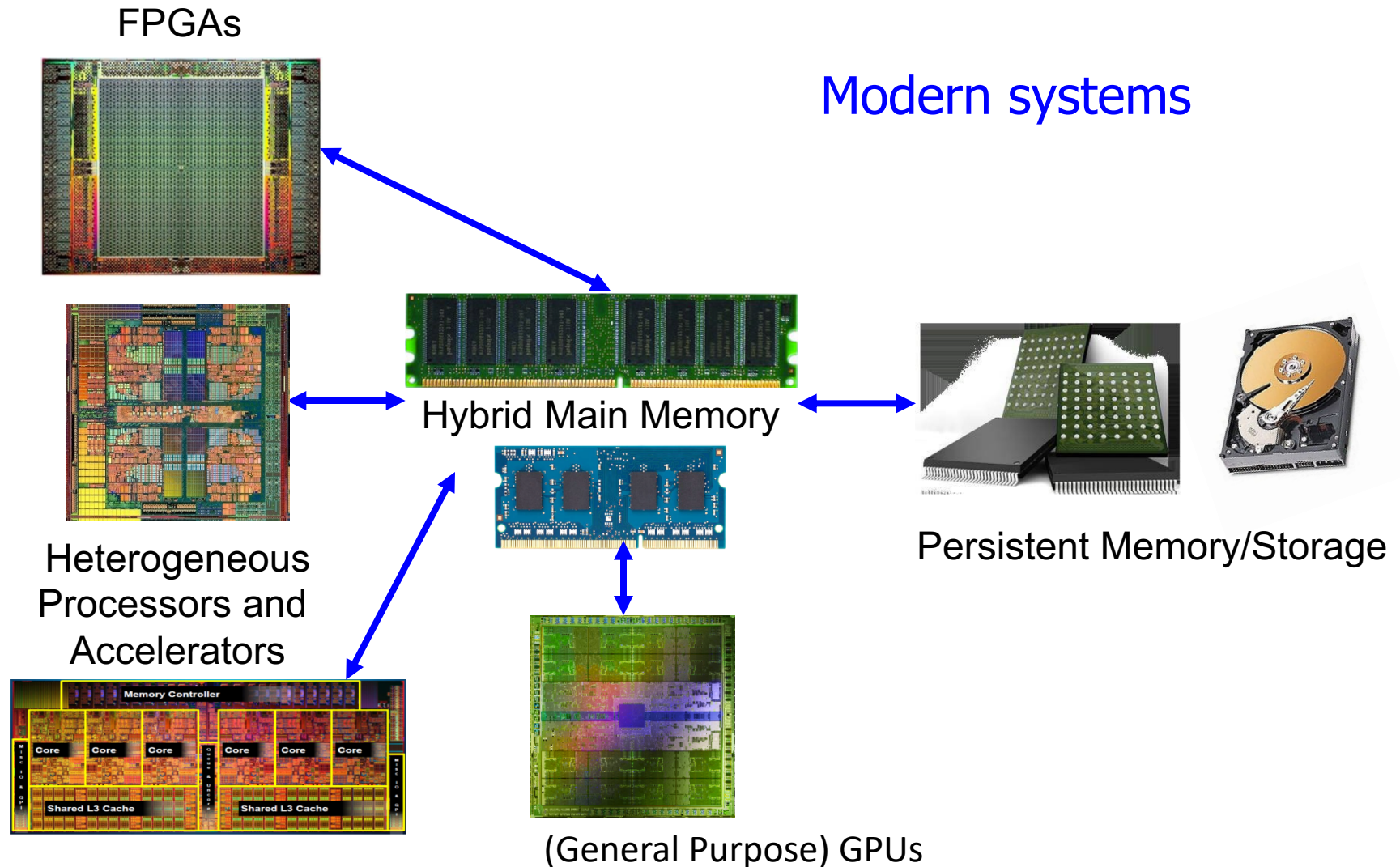
Main Memory



Storage (SSD/HDD)

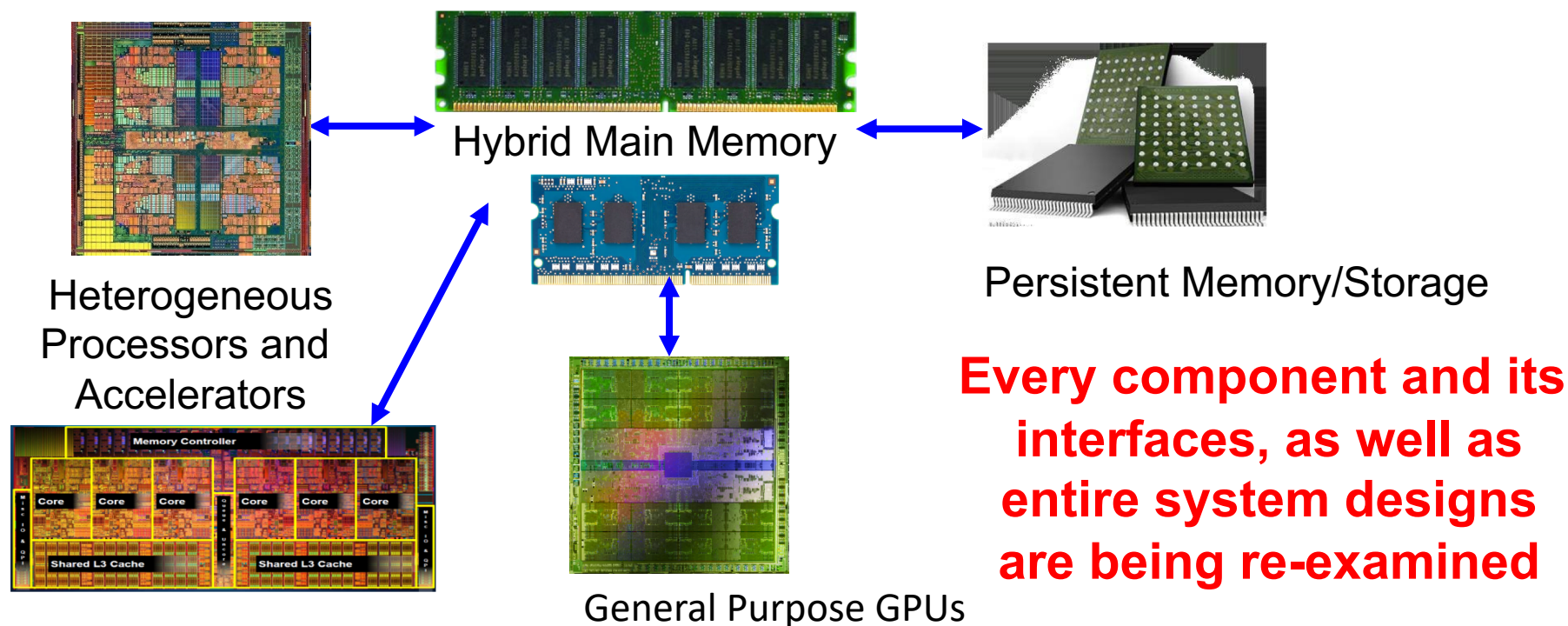


# Increasingly Complex Systems



# Computer Architecture Today

- Computing landscape is very different from 10-20 years ago
- Applications and technology both demand novel architectures



# Computer Architecture Today (II)

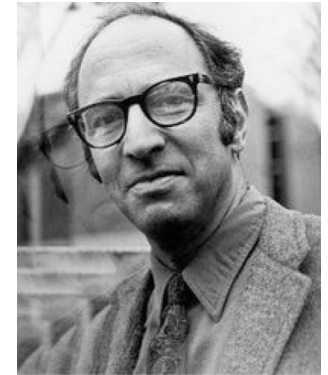
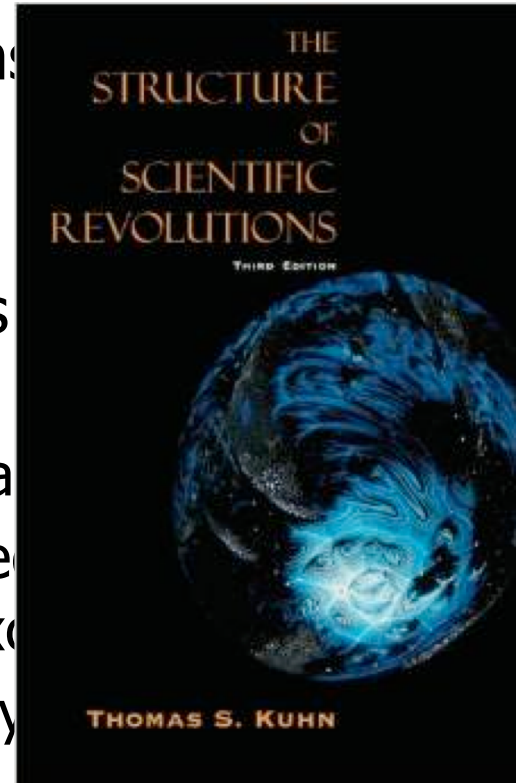
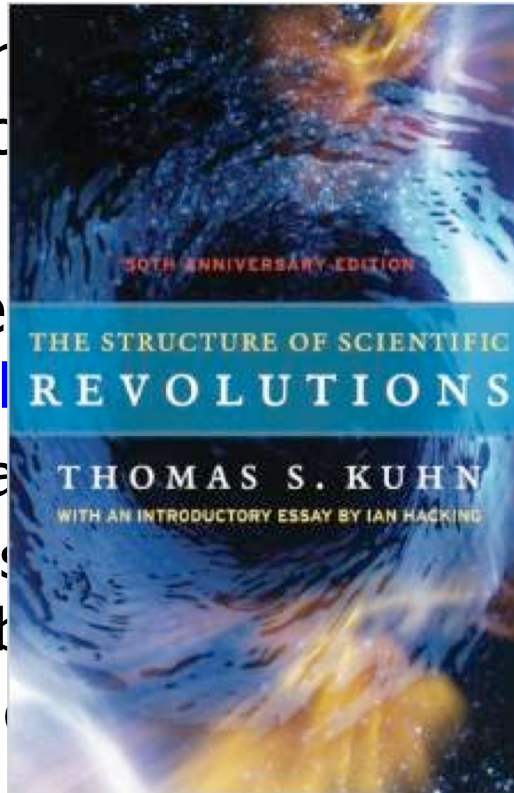
---

- You can revolutionize the way computers are built, if you understand both the hardware and the software (and change each accordingly)
- You can invent new paradigms for computation, communication, and storage
- Recommended book: Thomas Kuhn, “[The Structure of Scientific Revolutions](#)” (1962)
  - Pre-paradigm science: no clear consensus in the field
  - Normal science: dominant theory used to explain/improve things (business as usual); exceptions considered anomalies
  - Revolutionary science: underlying assumptions re-examined



# Computer Architecture Today (II)

- You can revolutionize the way computers are built, if you understand both the hardware and the software (and change each accordingly)
- You can improve communication
- Recommended reading:
  - **Scientific Revolutions** by Thomas S. Kuhn
    - Pre-para
    - Normal s
    - things (b
    - Revolution



ure of  
eld  
improve  
anomalies  
examined



# Takeaways

---

- It is an exciting time to be understanding and designing computing architectures
- Many challenging and exciting problems in platform design
  - That no one has tackled (or thought about) before
  - That can have huge impact on the world's future
- Driven by huge hunger for data (Big Data), new applications (ML/AI, graph analytics, genomics), ever-greater realism, ...
  - We can easily collect more data than we can analyze/understand
- Driven by significant difficulties in keeping up with that hunger at the technology layer
  - Five walls: Energy, reliability, complexity, security, scalability

# The Role of This Course

# Seminar in Comp Arch

---

- We will cover **fundamental** and **cutting-edge** research papers in computer architecture
- Multiple components that are aimed at improving students'
  - **technical skills** in computer architecture
  - **critical thinking and analysis**
  - **technical presentation** of concepts and papers
    - in both spoken and written forms
  - **familiarity with key research directions**

# Key Goal

---

(Learn how to)  
rigorously  
**analyze, present, discuss**  
papers and ideas  
in computer architecture

# Steps to Achieve the Key Goal

---

## ■ Steps for the Presenter

- ❑ Read
- ❑ Absorb, read more (other related works)
- ❑ Critically analyze; think; synthesize
- ❑ Prepare a clear and rigorous talk
- ❑ Present
- ❑ Answer questions
- ❑ Analyze and synthesize (in meeting, after, and at course end)

## ■ Steps for the Participants

- ❑ Discuss
- ❑ Ask questions
- ❑ Analyze and synthesize (in meeting, after, and at course end)

# Topics of Papers and Discussion

---

- hardware security;
- architectural acceleration mechanisms for key applications like machine learning, graph processing and bioinformatics;
- memory systems;
- interconnects;
- processing inside memory;
- various fundamental and emerging ideas/paradigms in computer architecture;
- hardware/software co-design and cooperation;
- fault tolerance;
- energy efficiency;
- heterogeneous and parallel systems;
- new execution models, etc.



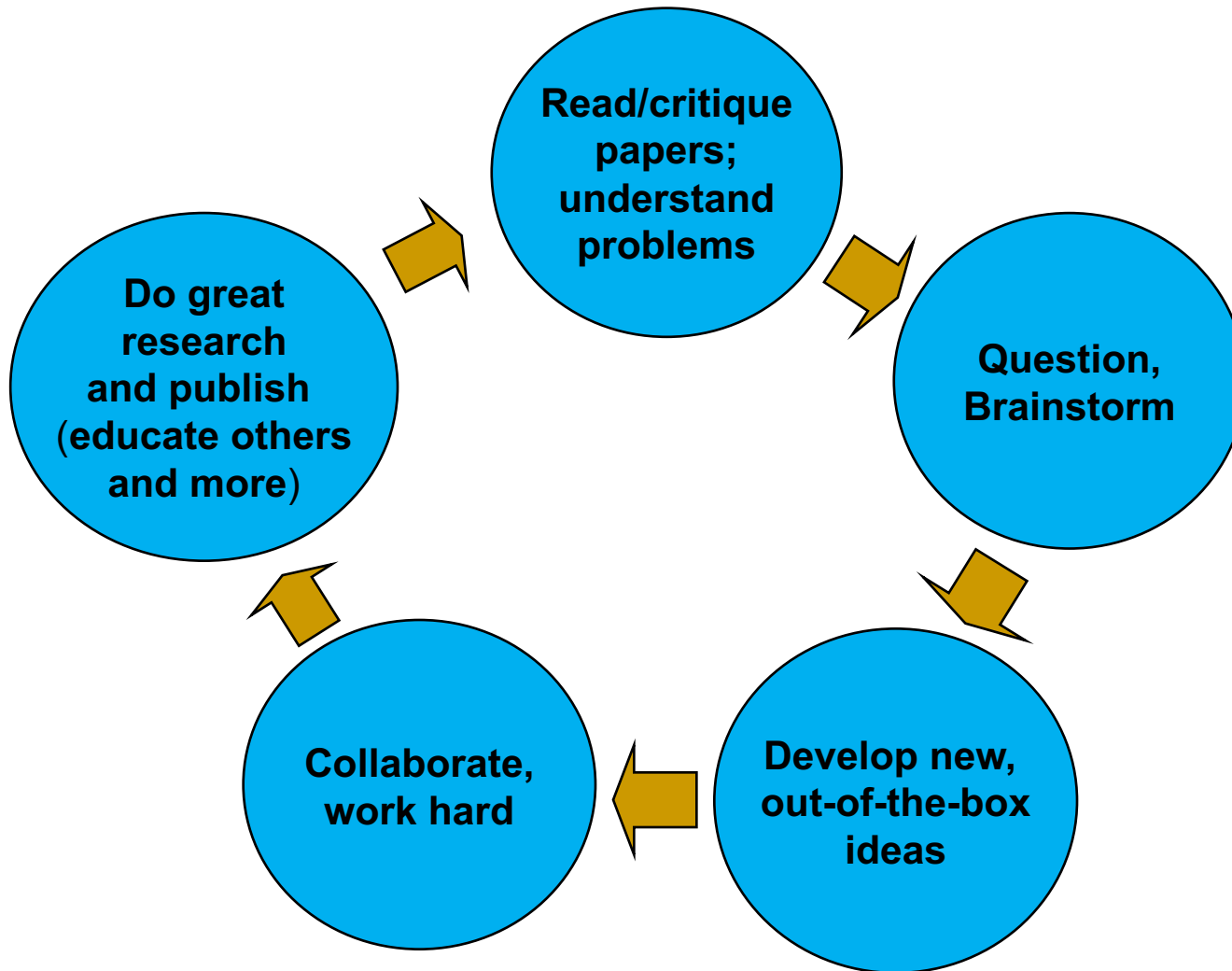
# Recap: Some Goals of This Course

---

- Teach/enable/empower you to:
  - ❑ Think critically
  - ❑ Think broadly
  - ❑ Learn how to understand, analyze and present papers and ideas
  - ❑ Get familiar with key first steps in research
  - ❑ Get familiar with key research directions

# The Virtuous Cycle of Scientific Progress

---



# Course Info and Logistics

# Course Info: Who Are We?

---



## ■ Onur Mutlu

- ❑ Full Professor @ ETH Zurich ITET (INFK), since September 2015
- ❑ Strecker Professor @ Carnegie Mellon University ECE/CS, 2009-2016, 2016-...
- ❑ PhD from UT-Austin, worked at Google, VMware, Microsoft Research, Intel, AMD
- ❑ <https://people.inf.ethz.ch/omutlu/>
- ❑ [omutlu@gmail.com](mailto:omutlu@gmail.com) (Best way to reach me)
- ❑ <https://people.inf.ethz.ch/omutlu/projects.htm>

## ■ Research and Teaching in:

- ❑ Computer architecture, computer systems, hardware security, bioinformatics
- ❑ Memory and storage systems
- ❑ Hardware security, safety, predictability
- ❑ Fault tolerance
- ❑ Hardware/software cooperation
- ❑ Architectures for bioinformatics, health, medicine
- ❑ ...

# Course Info: Who Are We?

---

## ■ Instructors:

- Dr. Mohammed Alser
- Dr. Juan Gomez Luna

## ■ Teaching Assistants

- Dr. Jisung Park,
- Dr. Jawad Haj-Yahya,
- Dr. Lois Orosa,
- Haiyu Mao
- Rahul Bera,
- João Dinis Ferreira,
- Geraldo Francisco,
- Can Firtina,
- Hasan Hassan,
- Konstantinos Kanellopoulos,
- Jeremie Kim,
- Nika Mansouri
- Minesh Patel,
- Gagandeep Singh,
- Kosta Stojiljkovic,
- Giray Yaglikci

## ■ Get to know them and their research

## ■ They will be your mentors – you will have to meet them at least twice before your presentations

# Course Requirements and Expectations

---

- Attendance required for all meetings
- Each student presents one paper
  - Prepare for presentation with engagement from the mentor
  - Full presentation + questions + discussion
- Non-presenters participate during the meeting
  - Ask questions, contribute thoughts/ideas
  - Better if you read/skim the paper beforehand
- Non-presenters take an online short quiz after each session
  - 5 MCQs for each presentation (2 hours to submit)
- Everyone comments on papers in the online review system
  - After presentation
- Write synthesis report at the end of semester
  - (sample synthesis report online)



# Course Website

---

- [https://safari.ethz.ch/architecture\\_seminar/fall2020](https://safari.ethz.ch/architecture_seminar/fall2020)
- All course materials to be posted
- Plus other useful information for the course
- Check frequently for announcements and due dates

# Homework 0

---

- Due September 24
  - [https://safari.ethz.ch/architecture\\_seminar/fall2020](https://safari.ethz.ch/architecture_seminar/fall2020)
- Information about yourself
- All future grading is predicated on homework 0
- If it is not submitted on time, we cannot schedule you for a presentation.

# Paper Review Preferences

---

- Due September 24
- Check the website for instructions
- If it is not submitted on time, we cannot schedule you for a presentation.

# How to Deliver a Good Talk

# Anatomy of a Good Paper Review (Talk)

---

- 0: Title, Authors, Venue
- 1: Summary
  - What is the problem the paper is trying to solve?
  - What are the key ideas of the paper? Key insights?
  - What are the key mechanisms? What is the implementation?
  - What are the key results? Key conclusions?
- 2: Strengths (most important ones)
  - Does the paper solve the problem well? Is it well written? ...
- 3: Weaknesses (most important ones)
  - This is where you should **think critically**. Every paper/idea has a weakness. This does not mean the paper is necessarily bad. It means there is room for improvement and future research can accomplish this.
- 4: Thoughts/Ideas: Can you do better? Present your ideas.
- 5: Takeaways: What you learned/enjoyed/disliked? Why?
- 6: Discussion starters and questions.
- Review should be short and concise (20 minutes or < one page)

# Suggested Paper Discussion Format

---

- Problem & Goal
- Key Ideas/solution
- Novelty
- Mechanisms & Implementation
- Major Results
- Takeaways/Conclusions

**~20-25 minute  
Summary**

- Strengths
- Weaknesses
- Alternatives
- New ideas/problems
- Brainstorming and Discussion

**~10 min Critique  
plus  
~15 min Discussion**

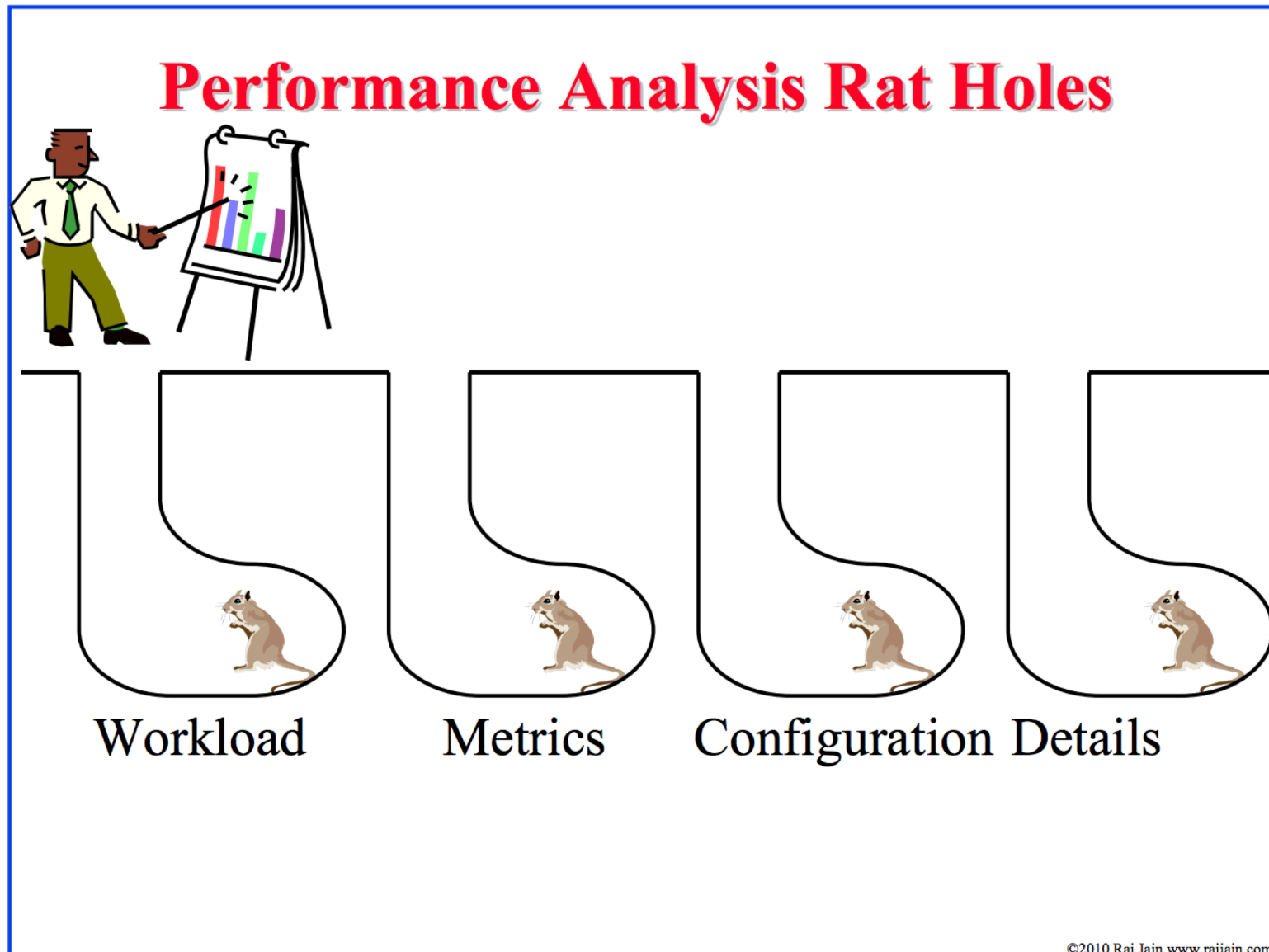


# More Advice on Paper Review Talk

---

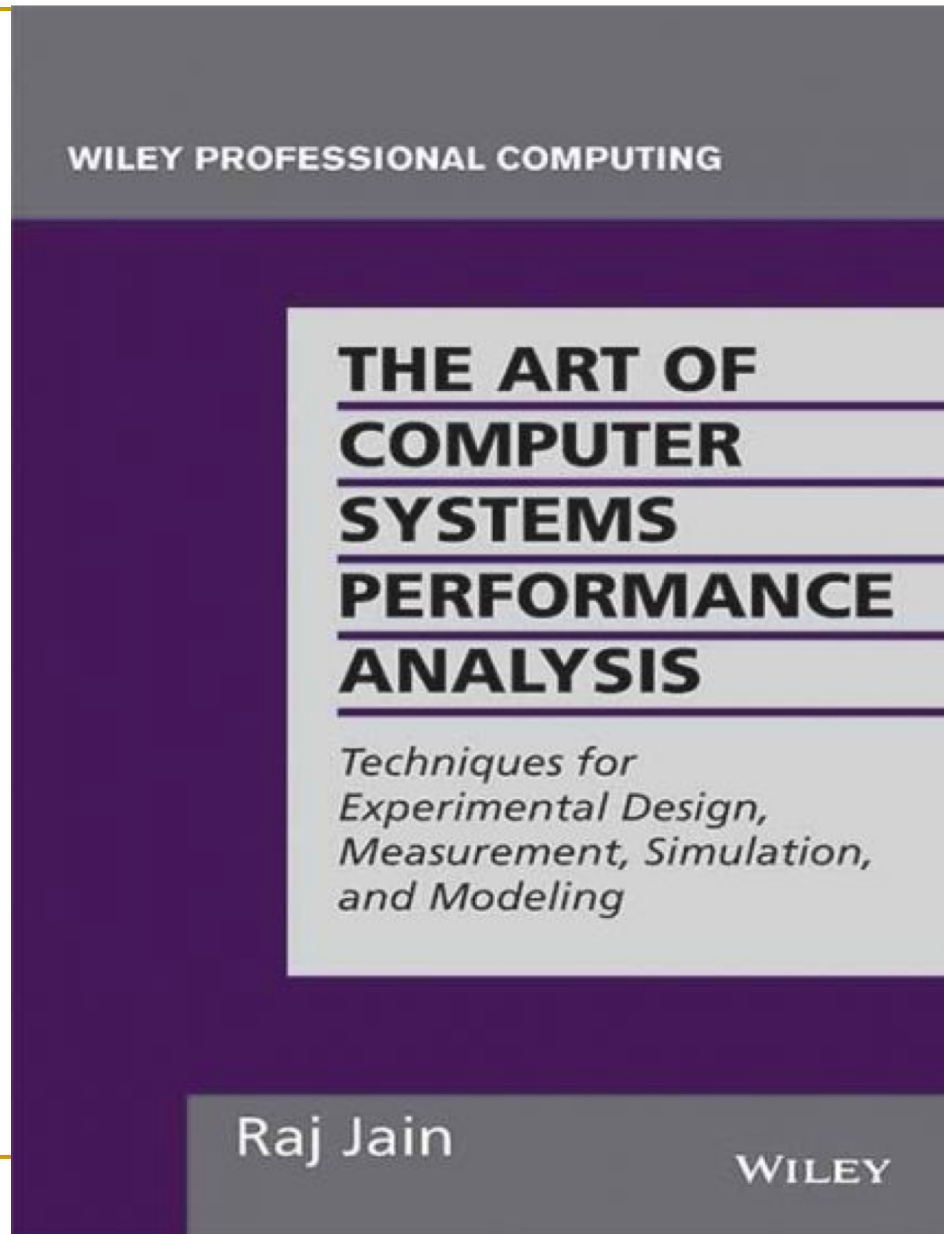
- When doing the paper reviews and analyses, be very critical
- Always think about better ways of solving the problem or related problems
  - Question the problem as well
  - Read background papers (both past and future)
- This is how things progress in science and engineering (or anywhere), and how you can make big leaps
  - By critical analysis
- A few sample text reviews provided online

# Try to Avoid Rat Hole Discussions



# Aside: A Recommended Book

---



Raj Jain, "[The Art of Computer Systems Performance Analysis](#)," Wiley, 1991.

## 10.8 DECISION MAKER'S GAMES

Even if the performance analysis is correctly done and presented, it may not be enough to persuade your audience—the decision makers—to follow your recommendations. The list shown in Box 10.2 is a compilation of reasons for rejection heard at various performance analysis presentations. You can use the list by presenting it immediately and pointing out that the reason for rejection is not new and that the analysis deserves more consideration. Also, the list is helpful in getting the competing proposals rejected!

There is no clear end of an analysis. Any analysis can be rejected simply on the grounds that the problem needs more analysis. This is the first reason listed in Box 10.2. The second most common reason for rejection of an analysis and for endless debate is the workload. Since workloads are always based on the past measurements, their applicability to the current or future environment can always be questioned. Actually workload is one of the four areas of discussion that lead a performance presentation into an endless debate. These “rat holes” and their relative sizes in terms of time consumed are shown in Figure 10.26. Presenting this cartoon at the beginning of a presentation helps to avoid these areas.

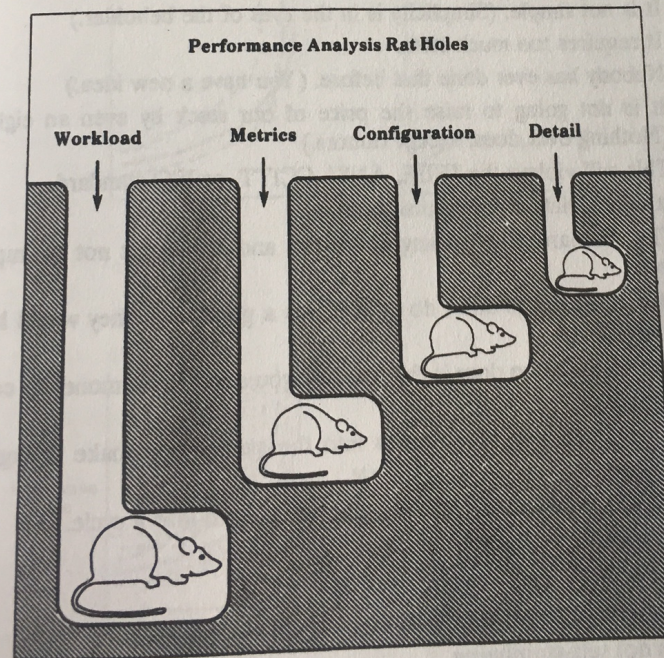


FIGURE 10.26 Four issues in performance presentations that commonly lead to endless discussion.

Raj Jain, "The Art of Computer Systems Performance Analysis," Wiley, 1991.



**Box 10.2 Reasons for Not Accepting the Results of an Analysis**

1. This needs more analysis.
2. You need a better understanding of the workload.
3. It improves performance only for long I/O's, packets, jobs, and files, and most of the I/O's, packets, jobs, and files are short.
4. It improves performance only for short I/O's, packets, jobs, and files, but who cares for the performance of short I/O's, packets, jobs, and files; its the long ones that impact the system.
5. It needs too much memory/CPU/bandwidth and memory/CPU/bandwidth isn't free.
6. It only saves us memory/CPU/bandwidth and memory/CPU/bandwidth is cheap.
7. There is no point in making the networks (similarly, CPUs/disks/...) faster; our CPUs/disks (any component other than the one being discussed) aren't fast enough to use them.
8. It improves the performance by a factor of  $x$ , but it doesn't really matter at the user level because everything else is so slow.
9. It is going to increase the complexity and cost.
10. Let us keep it simple stupid (and your idea is not stupid).
11. It is not simple. (Simplicity is in the eyes of the beholder.)
12. It requires too much state.
13. Nobody has ever done that before. (You have a new idea.)
14. It is not going to raise the price of our stock by even an eighth. (Nothing ever does, except rumors.)
15. This will violate the IEEE, ANSI, CCITT, or ISO standard.
16. It may violate some future standard.
17. The standard says nothing about this and so it must not be important.
18. Our competitors don't do it. If it was a good idea, they would have done it.
19. Our competition does it this way and you don't make money by copying others.
20. It will introduce randomness into the system and make debugging difficult.
21. It is too deterministic; it may lead the system into a cycle.
22. It's not interoperable.
23. This impacts hardware.
24. That's beyond today's technology.
25. It is not self-stabilizing.
26. Why change—it's working OK.

Raj Jain, "The Art of Computer Systems Performance Analysis," Wiley, 1991.

# More Advice on Talks

---

- Kayvon Fatahalian, “Tips for Giving Clear Talks”
  - <http://graphics.stanford.edu/~kayvonf/misc/cleartalktips.pdf>
  - Many useful and simple principles here

**“Every sentence matters”**

**“The audience prefers not to think” (about things you can just tell them)**

**“Surprises are bad”: say why before what  
(indicate why you are saying something before you say it)**

**Explain every figure, graph, or equation**

**When improving the talk, the audience is always right**

# Who Painted This Painting?

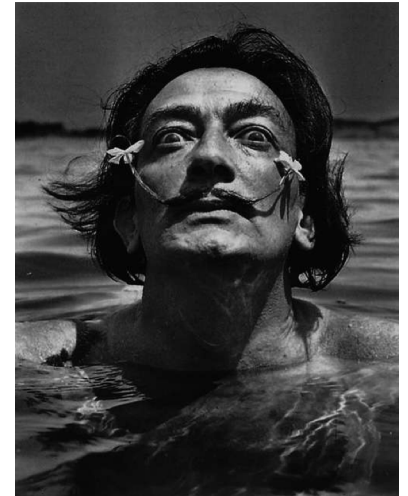
---



Salvador Dali @ 1924



# What About This?



Salvador Dali @ 1937

# Takeaway

---

Learn the basic principles

**before** you

*consciously* choose to break them

# How to Participate

# How to Make the Best Out of This?

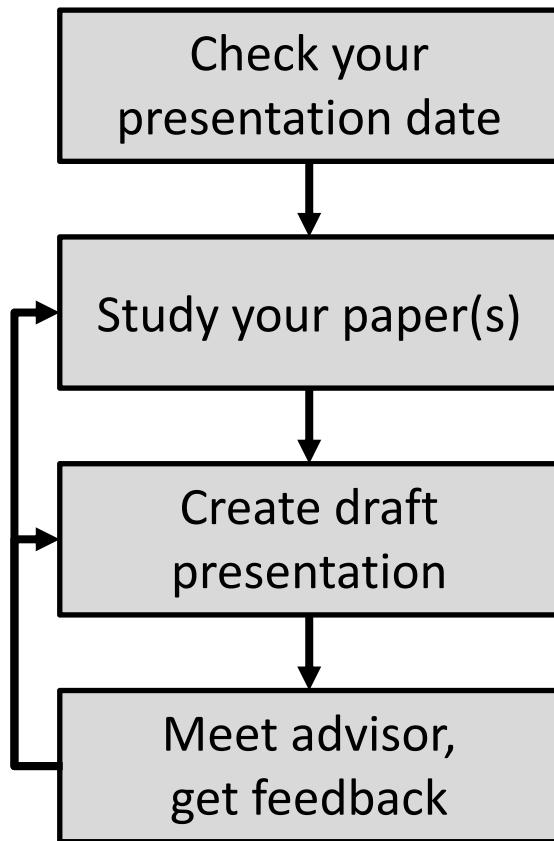
---

- Come prepared → Read and critically evaluate the paper
- Think new ideas
- Bring discussion points and questions; read other papers
- Be critical
- Brainstorm – be open to new ideas
- Pay attention and discuss+contribute
- Participate online before and after each meeting

# Guided Talk Preparation

# Preparing a Talk

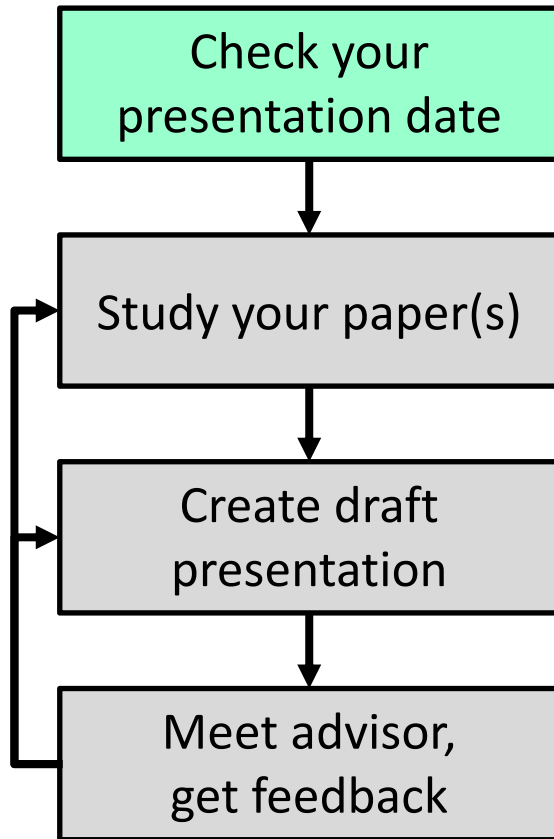
---





# Preparing a Talk: Start Early

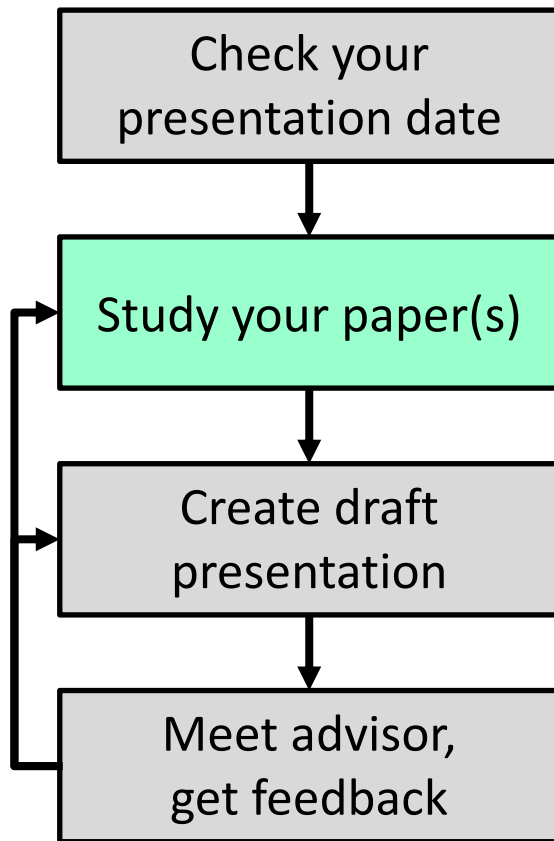
---



- Preparing a good presentation takes time
- Start early!

# Preparing a Talk: Study Paper

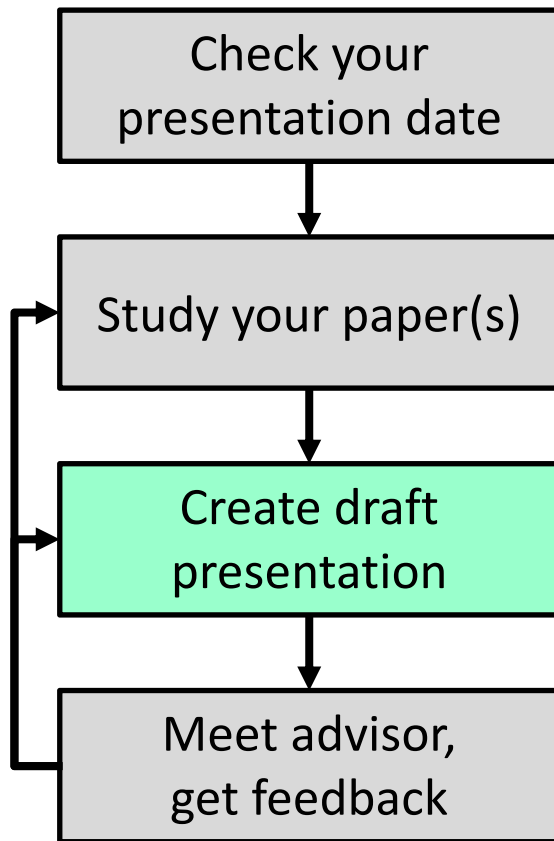
---



- 3 'C's of reading
  - ❑ *Carefully*: look up terms, possibly read cited papers
  - ❑ *Critically*: find limitations, flaws
  - ❑ *Creatively*: think of improvements
- Try examples by hand
- Try tools if available
- Consult with TA if questions

# Preparing a Talk: Create Draft

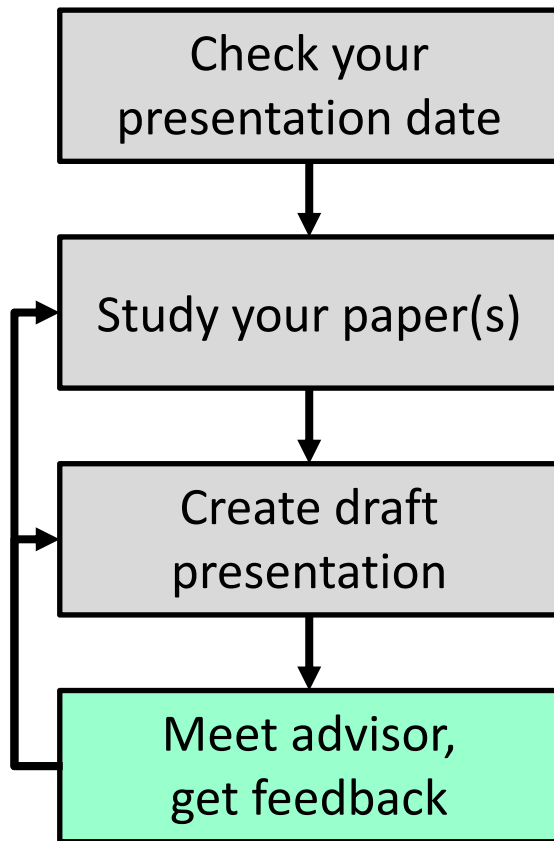
---



- Explain the motivation for the work
- Clearly present the technical solution and results
  - Include a demo if appropriate
- Outline limitations or improvements
- Focus on the key concepts
  - Do not present all of the details

# Preparing a Talk: Get Feedback

---



- Prepare for the meeting
  - Schedule early
  - Send slides in advance
  - Write down questions
- Make sure you address feedback
  - Take notes
- Meetings are mandatory!
  - At least one week before the talk
  - Two meetings

# Grading and Feedback

# Grading Rubric

---

- Quality of your presentation (60%)
  - ❑ How well did you understand the material?
  - ❑ How well did you present it?
  - ❑ How well did you answer the questions?
  - ❑ Be prepared to explain technical terms
  - ❑ **We will take into account** the difficulty of the paper and the time you had to prepare.
- Quality of the final synthesis paper (30%)
  - ❑ How well did you understand some of the papers presented during the seminar?
- Attendance & Quizzes (10%)
- Participation (during class and online) (BONUS 10%)
  - ❑ Did you ask good questions?
  - ❑ Did you attend all sessions?

# Feedback

---

- We will try to (briefly) discuss strengths/weaknesses of your talk in class
  - Let us know upfront if you would prefer **not** to
- You can arrange a meeting with your TA to get feedback



# Expected Schedule

# Schedule

---

- We will meet once a week, with two presentations per session
  - *Next meeting next week*
  - *Your presentations start on 8 October*
  - 22 presentations in total
  - Each presentation 50 minutes including questions and discussion
- Paper assignment
  - Will be done online
  - Study the list of papers
  - **Check your email** and be responsive

# Homework 0

---

- Due September 24
  - [https://safari.ethz.ch/architecture\\_seminar/fall2020](https://safari.ethz.ch/architecture_seminar/fall2020)
- Information about yourself
- All future grading is predicated on homework 0
- If it is not submitted on time, we cannot schedule you for a presentation.

# Paper Review Preferences

---

- Due September 24
- Check the website for instructions
- If it is not submitted on time, we cannot schedule you for a presentation.

# Seminar in Computer Architecture Meeting 1: Introduction

Prof. Onur Mutlu

ETH Zürich

Fall 2020

17 September 2020