

# Computer Architecture

## Lecture 4a: Memory Systems: Solution Directions

Prof. Onur Mutlu

ETH Zürich

Fall 2020

25 September 2020

# Solving the Memory Problem

# How Do We Solve The Memory Problem?

---

- **Fix it:** Make memory and controllers more intelligent
  - **New interfaces, functions, architectures:** system-mem codesign
- **Eliminate or minimize it:** Replace or (more likely) augment DRAM with a different technology
  - **New technologies and system-wide rethinking** of memory & storage
- **Embrace it:** Design heterogeneous memories (none of which are perfect) and map data intelligently across them
  - **New models for data management and maybe usage**
- ...

# How Do We Solve The Memory Problem?

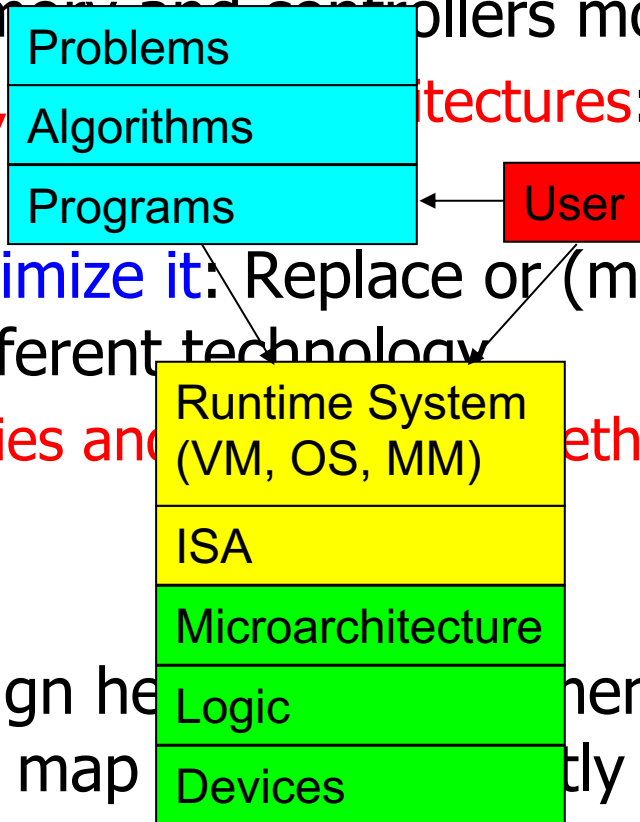
---

- **Fix it:** Make memory and controllers more intelligent
  - **New interfaces, functions, architectures:** system-mem codesign
- **Eliminate or minimize it:** Replace or (more likely) augment DRAM with a different technology
  - **New technologies and system-wide rethinking** of memory & storage
- **Embrace it:** Design heterogeneous memories (none of which are perfect) and map data intelligently across them
  - **New models for data management and maybe usage**

**Solutions (to memory scaling) require software/hardware/device cooperation**

# How Do We Solve The Memory Problem?

- **Fix it:** Make memory and controllers more intelligent
  - **New interfaces, architectures:** system-mem codesign
- **Eliminate or minimize it:** Replace or (more likely) augment DRAM with a different technology
  - **New technologies and storage**
- **Embrace it:** Design heterogeneous memories (none of which are perfect) and map applications directly across them
  - **New models for data management and maybe usage**



**Solutions (to memory scaling) require software/hardware/device cooperation**

# Solution 1: New Memory Architectures

---

- Overcome memory shortcomings with
  - ❑ Memory-centric system design
  - ❑ Novel memory architectures, interfaces, functions
  - ❑ Better waste management (efficient utilization)
- Key issues to tackle
  - ❑ Enable reliability at low cost → high capacity
  - ❑ Reduce energy
  - ❑ Reduce latency
  - ❑ Improve bandwidth
  - ❑ Reduce waste (capacity, bandwidth, latency)
  - ❑ Enable computation close to data

# Solution 1: New Memory Architectures

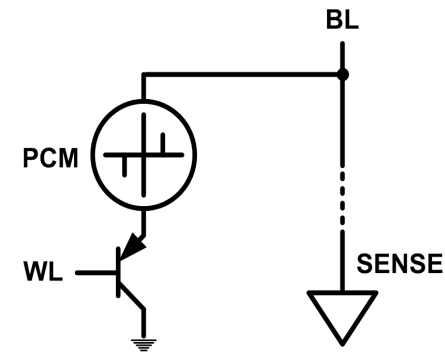
- Liu+, "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
- Kim+, "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012.
- Lee+, "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.
- Liu+, "An Experimental Study of Data Retention Behavior in Modern DRAM Devices," ISCA 2013.
- Seshadri+, "RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," MICRO 2013.
- Pekhimenko+, "Linearly Compressed Pages: A Main Memory Compression Framework," MICRO 2013.
- Chang+, "Improving DRAM Performance by Parallelizing Refreshes with Accesses," HPCA 2014.
- Khan+, "The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study," SIGMETRICS 2014.
- Luo+, "Characterizing Application Memory Error Vulnerability to Optimize Data Center Cost," DSN 2014.
- Kim+, "Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors," ISCA 2014.
- Lee+, "Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common-Case," HPCA 2015.
- Qureshi+, "AVATAR: A Variable-Retention-Time (VRT) Aware Refresh for DRAM Systems," DSN 2015.
- Meza+, "Revisiting Memory Errors in Large-Scale Production Data Centers: Analysis and Modeling of New Trends from the Field," DSN 2015.
- Kim+, "Ramulator: A Fast and Extensible DRAM Simulator," IEEE CAL 2015.
- Seshadri+, "Fast Bulk Bitwise AND and OR in DRAM," IEEE CAL 2015.
- Ahn+, "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing," ISCA 2015.
- Ahn+, "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture," ISCA 2015.
- Lee+, "Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM," PACT 2015.
- Seshadri+, "Gather-Scatter DRAM: In-DRAM Address Translation to Improve the Spatial Locality of Non-unit Strided Accesses," MICRO 2015.
- Lee+, "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost," TACO 2016.
- Hassan+, "ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality," HPCA 2016.
- Chang+, "Low-Cost Inter-Linked Subarrays (LISA): Enabling Fast Inter-Subarray Data Migration in DRAM," HPCA 2016.
- Chang+, "Understanding Latency Variation in Modern DRAM Chips: Experimental Characterization, Analysis, and Optimization," SIGMETRICS 2016.
- Khan+, "PARBOR: An Efficient System-Level Technique to Detect Data Dependent Failures in DRAM," DSN 2016.
- Hsieh+, "Transparent Offloading and Mapping (TOM): Enabling Programmer-Transparent Near-Data Processing in GPU Systems," ISCA 2016.
- Hashemi+, "Accelerating Dependent Cache Misses with an Enhanced Memory Controller," ISCA 2016.
- Boroumand+, "LazyPIM: An Efficient Cache Coherence Mechanism for Processing-in-Memory," IEEE CAL 2016.
- Pattnaik+, "Scheduling Techniques for GPU Architectures with Processing-In-Memory Capabilities," PACT 2016.
- Hsieh+, "Accelerating Pointer Chasing in 3D-Stacked Memory: Challenges, Mechanisms, Evaluation," ICCD 2016.
- Hashemi+, "Continuous Runahead: Transparent Hardware Acceleration for Memory Intensive Workloads," MICRO 2016.
- Khan+, "A Case for Memory Content-Based Detection and Mitigation of Data-Dependent Failures in DRAM," IEEE CAL 2016.
- Hassan+, "SoftMC: A Flexible and Practical Open-Source Infrastructure for Enabling Experimental DRAM Studies," HPCA 2017.
- Mutlu+, "The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser," DATE 2017.
- Lee+, "Design-Induced Latency Variation in Modern DRAM Chips: Characterization, Analysis, and Latency Reduction Mechanisms," SIGMETRICS 2017.
- Chang+, "Understanding Reduced-Voltage Operation in Modern DRAM Devices: Experimental Characterization, Analysis, and Mechanisms," SIGMETRICS 2017.
- Patel+, "The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions," ISCA 2017.
- Seshadri and Mutlu, "Simple Operations in Memory to Reduce Data Movement," ADCOM 2017.
- Liu+, "Concurrent Data Structures for Near-Memory Computing," SPAA 2017.
- Khan+, "Detecting and Mitigating Data-Dependent DRAM Failures by Exploiting Current Memory Content," MICRO 2017.
- Seshadri+, "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," MICRO 2017.
- Kim+, "GRIM-Filter: Fast Seed Location Filtering in DNA Read Mapping Using Processing-in-Memory Technologies," BMC Genomics 2018.
- Kim+, "The DRAM Latency PUF: Quickly Evaluating Physical Unclonable Functions by Exploiting the Latency-Reliability Tradeoff in Modern DRAM Devices," HPCA 2018.
- Boroumand+, "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks," ASPLOS 2018.
- Das+, "VRL-DRAM: Improving DRAM Performance via Variable Refresh Latency," DAC 2018.
- Ghose+, "What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study," SIGMETRICS 2018.
- Kim+, "Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines," ICCD 2018.
- Wang+, "Reducing DRAM Latency via Charge-Level-Aware Look-Ahead Partial Restoration," MICRO 2018.
- Kim+, "D-RaNGe: Using Commodity DRAM Devices to Generate True Random Numbers with Low Latency and High Throughput," HPCA 2019.
- Singh+, "NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning," DAC 2019.
- Ghose+, "Demystifying Workload-DRAM Interactions: An Experimental Study," SIGMETRICS 2019.
- Patel+, "Understanding and Modeling On-Die Error Correction in Modern DRAM: An Experimental Study Using Real Devices," DSN 2019.
- Boroumand+, "CoNDA: Efficient Cache Coherence Support for Near-Data Accelerators," ISCA 2019.
- Hassan+, "CROW: A Low-Cost Substrate for Improving DRAM Performance, Energy Efficiency, and Reliability," ISCA 2019.
- Mutlu and Kim, "RowHammer: A Retrospective," TCAD 2019.
- Mutlu+, "Processing Data Where It Makes Sense: Enabling In-Memory Computation," MICRO 2019.
- Seshadri and Mutlu, "In-DRAM Bulk Bitwise Execution Engine," ADCOM 2020.
- Koppula+, "EDEN: Energy-Efficient, High-Performance Neural Network Inference Using Approximate DRAM," MICRO 2019.
- Rezaei+, "NoM: Network-on-Memory for Inter-Bank Data Transfer in Highly-Banked Memories," CAL 2020.
- Frigo+, "TRRespass: Exploiting the Many Sides of Target Row Refresh," S&P 2020.
- Coljocar+, "Are We Susceptible to Rowhammer? An End-to-End Methodology for Cloud Providers," S&P 2020.
- Luo+, "CLR-DRAM: A Low-Cost DRAM Architecture Enabling Dynamic Capacity-Latency Trade-Off," ISCA 2020.
- Kim+, "Revisiting RowHammer: An Experimental Analysis of Modern Devices and Mitigation Techniques," ISCA 2020.
- Wang+, "FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching," MICRO 2020.
- Patel+, "Bit-Exact ECC Recovery (BEER): Determining DRAM On-Die ECC Functions by Exploiting DRAM Data Retention Characteristics," MICRO 2020.
- Avoid DRAM:
  - Seshadri+, "The Evicted-Address Filter: A Unified Mechanism to Address Both Cache Pollution and Thrashing," PACT 2012.
  - Pekhimenko+, "Base-Delta-Immediate Compression: Practical Data Compression for On-Chip Caches," PACT 2012.
  - Seshadri+, "The Dirty-Block Index," ISCA 2014.
  - Pekhimenko+, "Exploiting Compressed Block Size as an Indicator of Future Reuse," HPCA 2015.
  - Vijaykumar+, "A Case for Core-Assisted Bottleneck Acceleration in GPUs: Enabling Flexible Data Compression with Assist Warps," ISCA 2015.
  - Pekhimenko+, "Toggle-Aware Bandwidth Compression for GPUs," HPCA 2016.

# Solution 2: Emerging Memory Technologies

- Some emerging **resistive** memory technologies seem more scalable than DRAM (and they are non-volatile)

- Example: Phase Change Memory

- Data stored by changing phase of material
- Data read by detecting material's resistance
- Expected to scale to 9nm (2022 [ITRS 2009])
- Prototyped at 20nm (Raoux+, IBM JRD 2008)
- Expected to be denser than DRAM: can store multiple bits/cell



- But, emerging technologies have (many) shortcomings
  - Can they be enabled to replace/augment/surpass DRAM?



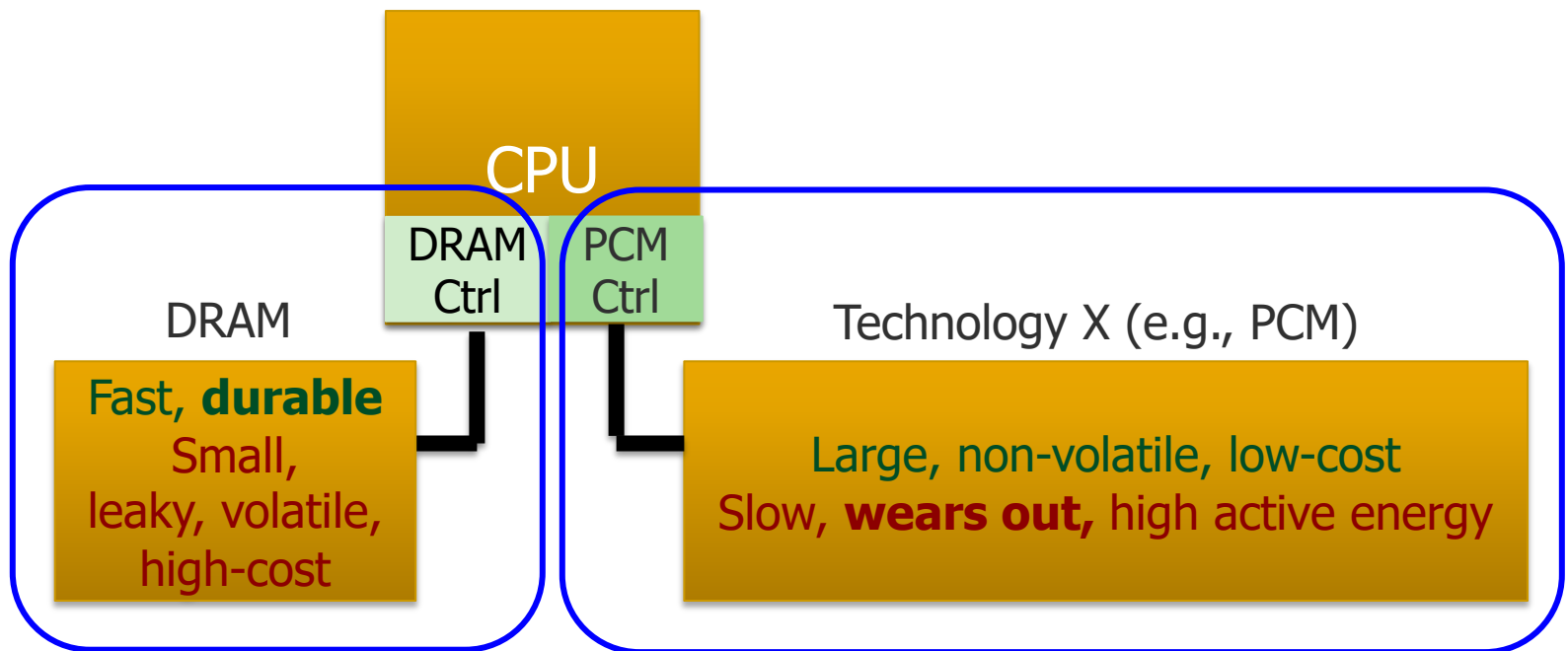
# Solution 2: Emerging Memory Technologies

---

- Lee+, “[Architecting Phase Change Memory as a Scalable DRAM Alternative](#),” ISCA’09, CACM’10, IEEE Micro’10.
- Meza+, “[Enabling Efficient and Scalable Hybrid Memories](#),” IEEE Comp. Arch. Letters 2012.
- Yoon, Meza+, “[Row Buffer Locality Aware Caching Policies for Hybrid Memories](#),” ICCD 2012.
- Kultursay+, “[Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative](#),” ISPASS 2013.
- Meza+, “[A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory](#),” WEED 2013.
- Lu+, “[Loose Ordering Consistency for Persistent Memory](#),” ICCD 2014.
- Zhao+, “[FIRM: Fair and High-Performance Memory Control for Persistent Memory Systems](#),” MICRO 2014.
- Yoon, Meza+, “[Efficient Data Mapping and Buffering Techniques for Multi-Level Cell Phase-Change Memories](#),” TACO 2014.
- Ren+, “[ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems](#),” MICRO 2015.
- Chauhan+, “[NVMove: Helping Programmers Move to Byte-Based Persistence](#),” INFLOW 2016.
- Li+, “[Utility-Based Hybrid Memory Management](#),” CLUSTER 2017.
- Yu+, “[Banshee: Bandwidth-Efficient DRAM Caching via Software/Hardware Cooperation](#),” MICRO 2017.
- Tavakkol+, “[MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices](#),” FAST 2018.
- Tavakkol+, “[FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives](#),” ISCA 2018.
- Sadrosadati+. “[LTRF: Enabling High-Capacity Register Files for GPUs via Hardware/Software Cooperative Register Prefetching](#),” ASPLOS 2018.
- Salkhordeh+, “[An Analytical Model for Performance and Lifetime Estimation of Hybrid DRAM-NVM Main Memories](#),” TC 2019.
- Wang+, “[Panthera: Holistic Memory Management for Big Data Processing over Hybrid Memories](#),” PLDI 2019.
- Song+, “[Enabling and Exploiting Partition-Level Parallelism \(PALP\) in Phase Change Memories](#),” CASES 2019.
- Liu+, “[Binary Star: Coordinated Reliability in Heterogeneous Memory Systems for High Performance and Scalability](#),” MICRO’19.
- Song+, “[Improving Phase Change Memory Performance with Data Content Aware Access](#),” ISMM 2020.

# Combination: Hybrid Memory Systems

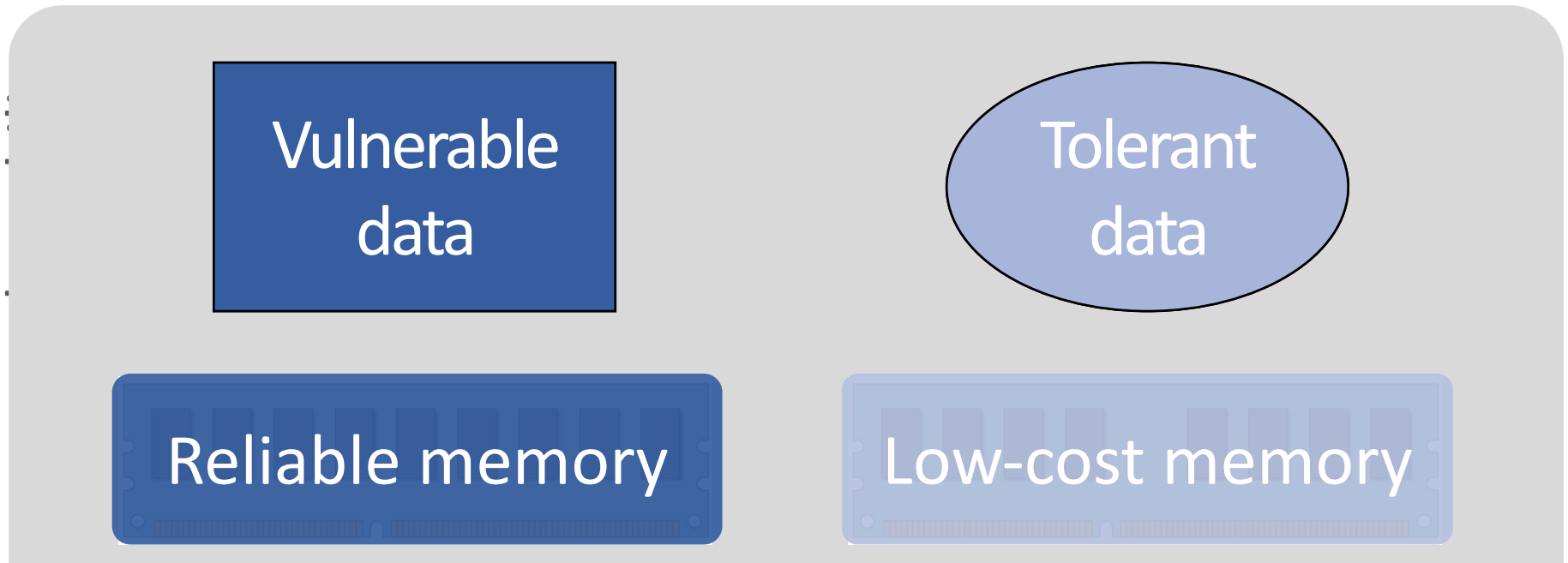
---



Hardware/software manage data allocation and movement  
to achieve the best of multiple technologies

Meza+, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters, 2012.  
Yoon, Meza et al., "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012 Best Paper Award.

# Exploiting Memory Error Tolerance with Hybrid Memory Systems



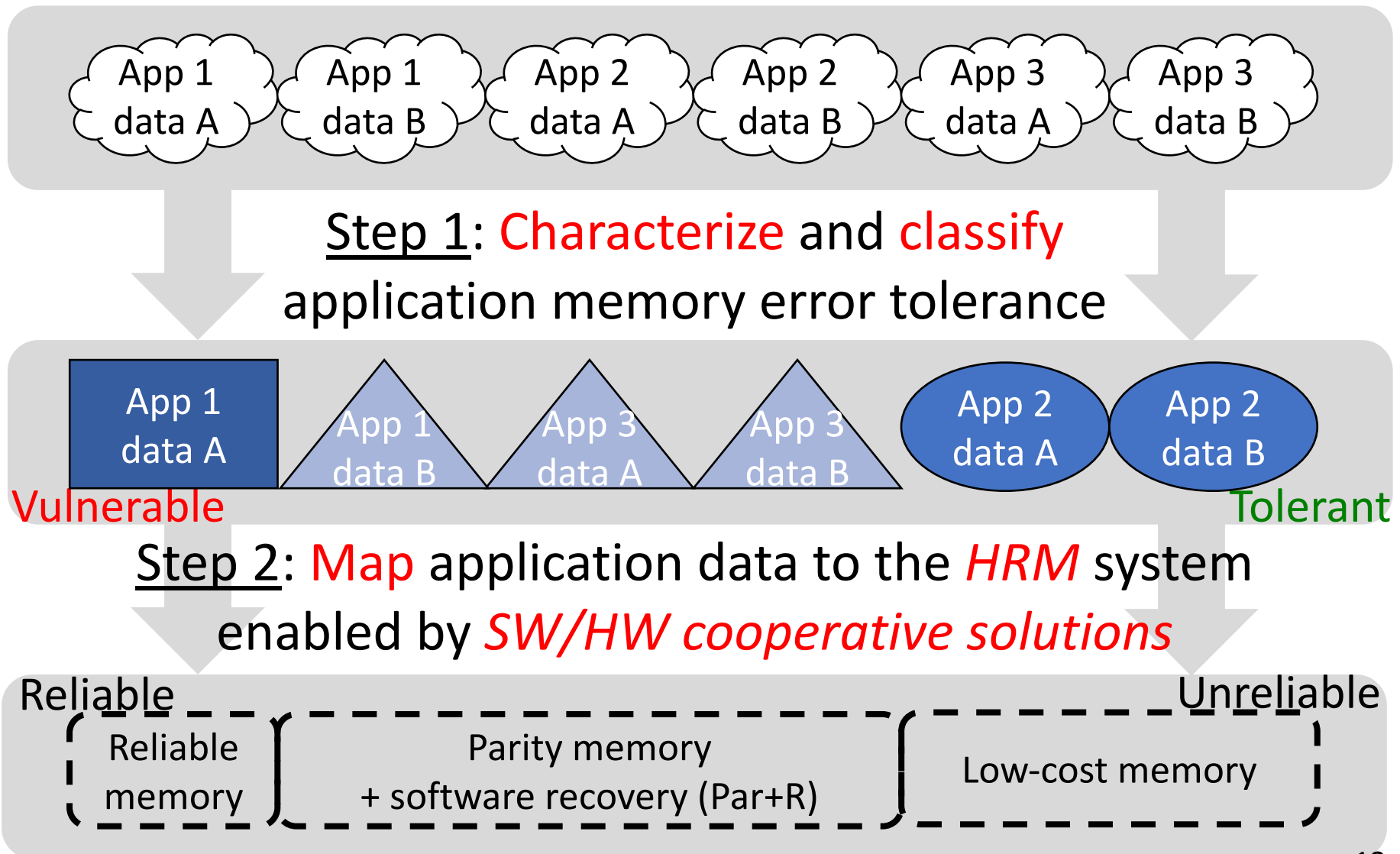
On Microsoft's Web Search workload

Reduces server hardware **cost** by **4.7 %**

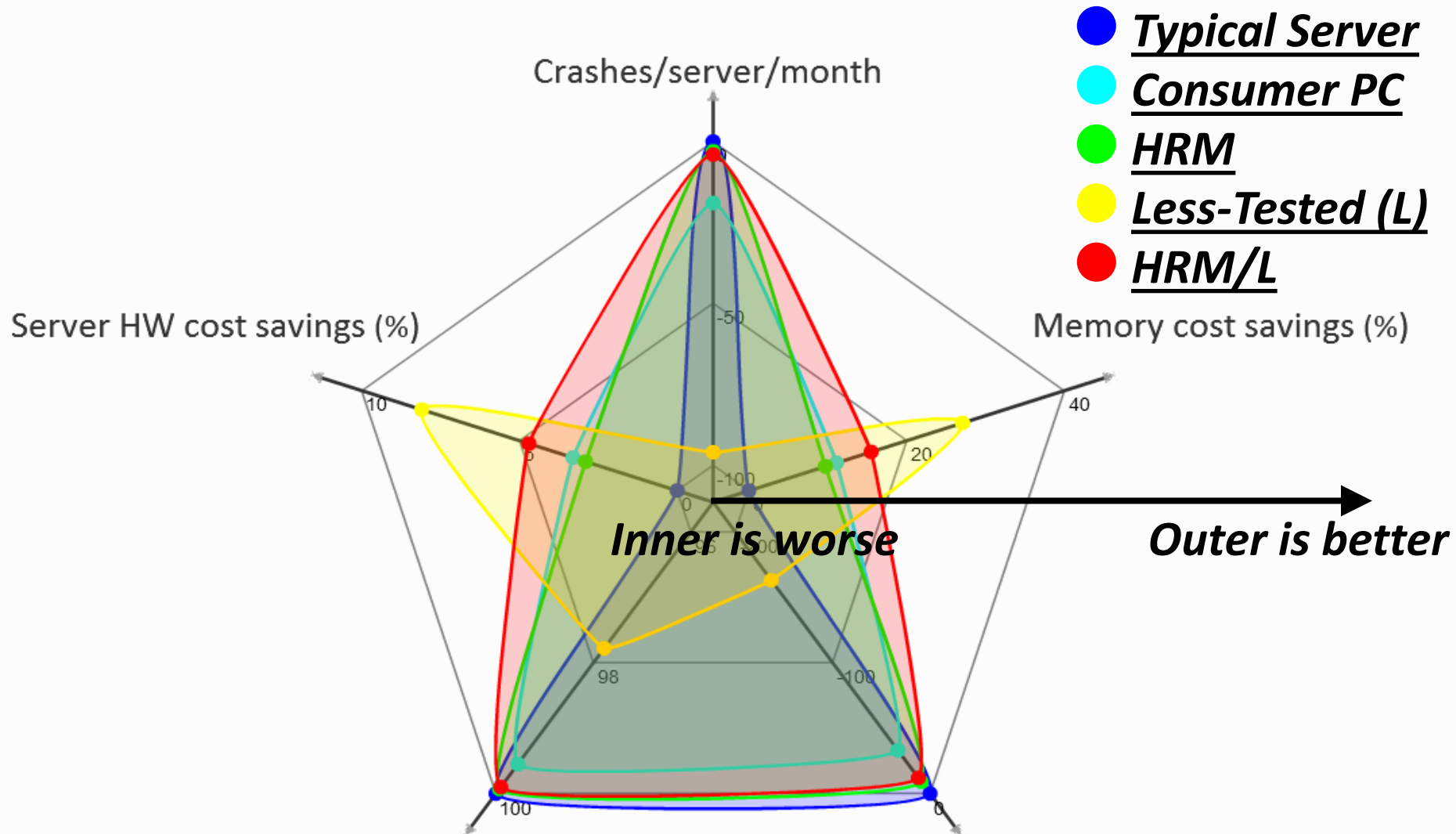
Achieves single server **availability** target of **99.90 %**

**Heterogeneous-Reliability Memory** [DSN 2014]

# Heterogeneous-Reliability Memory



# Evaluation Results



● ● Bigger area means better tradeoff

# More on 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

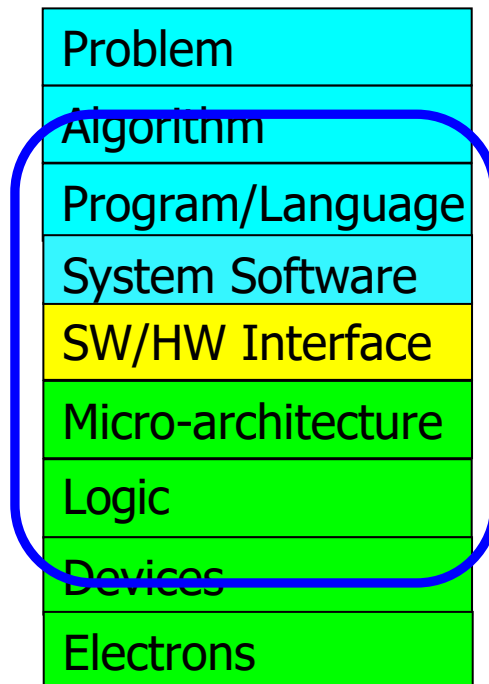
\*Microsoft Corporation, {srgovin, bsharma, marksan, kansal, jie.liu, bk Hessib, kvaid}@microsoft.com

# HRM is an Example of Our 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**

# An Orthogonal Issue: Memory Interference

---

- Problem: **Memory interference between cores is uncontrolled**
    - unfairness, starvation, low performance
    - **uncontrollable, unpredictable, vulnerable system**
  - Solution: **QoS-Aware Memory Systems**
    - Hardware designed to provide a configurable fairness substrate
      - Application-aware memory scheduling, partitioning, throttling
    - Software designed to configure the resources to satisfy different QoS goals
  - QoS-aware memory systems can provide predictable performance and higher efficiency
-



# Strong Memory Service Guarantees

---

- Goal: Satisfy performance/SLA requirements in the presence of shared main memory, heterogeneous agents, and hybrid memory/storage
- Approach:
  - Develop techniques/models to accurately estimate the performance loss of an application/agent in the presence of resource sharing
  - Develop mechanisms (hardware and software) to enable the resource partitioning/prioritization needed to achieve the required performance levels for all applications
  - All the while providing high system performance
- Subramanian et al., “MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems,” HPCA 2013.
- Subramanian et al., “The Application Slowdown Model,” MICRO 2015.

# Memory Controllers

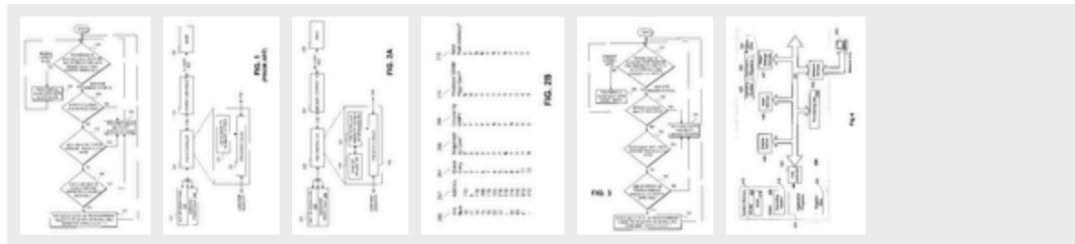
# It All Started with FSB Controllers (2001)

## Method and apparatus to control memory accesses

### Abstract

A method and apparatus for accessing memory comprising monitoring memory accesses from a hardware prefetcher and determining whether the memory accesses from the hardware prefetcher are used by an out-of-order core. A front side bus controller switches memory access modes from a minimize memory access latency mode to a maximize memory bus bandwidth mode if a percentage of the memory accesses generated by the hardware prefetcher are used by the out-of-order core.

### Images (6)



### Classifications

[G06F12/0215](#) Addressing or allocation; Relocation with look ahead addressing means

US6799257B2

United States

[Download PDF](#) [Find Prior Art](#) [Similar](#)

**Inventor:** [Eric A. Sprangle](#), [Onur Mutlu](#)

**Current Assignee :** [Intel Corp](#)

### Worldwide applications

2002 • [US](#) 2003 • [AU](#) [JP](#) [DE](#) [KR](#) [CN](#) [WO](#) [GB](#) [TW](#) 2004 • [US](#)  
2005 • [HK](#)

### Application US10/079,967 events ⓘ

2002-02-21 • Application filed by Intel Corp

2002-02-21 • Priority to US10/079,967

2002-04-25 • Assigned to INTEL CORPORATION ⓘ

# Memory Performance Attacks [USENIX SEC'07]

---

- Thomas Moscibroda and Onur Mutlu,  
**"Memory Performance Attacks: Denial of Memory Service in Multi-Core Systems"**  
*Proceedings of the 16th USENIX Security Symposium (**USENIX SECURITY**), pages 257-274, Boston, MA, August 2007. [Slides](#) ([ppt](#))*

## **Memory Performance Attacks: Denial of Memory Service in Multi-Core Systems**

*Thomas Moscibroda   Onur Mutlu  
Microsoft Research  
{moscitho,onur}@microsoft.com*

# STFM [MICRO'07]

---

- Onur Mutlu and Thomas Moscibroda,  
**"Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors"**  
*Proceedings of the 40th International Symposium on Microarchitecture (MICRO)*, pages 146-158, Chicago, IL, December 2007. [[Summary](#)] [[Slides \(ppt\)](#)]

## Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors

Onur Mutlu   Thomas Moscibroda

Microsoft Research  
{onur,moscitho}@microsoft.com

---

- Onur Mutlu and Thomas Moscibroda,  
**"Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems"**  
*Proceedings of the 35th International Symposium on Computer Architecture (ISCA)*, pages 63-74, Beijing, China, June 2008.  
[[Summary](#)] [[Slides \(ppt\)](#)]

## Parallelism-Aware Batch Scheduling:

## Enhancing both Performance and Fairness of Shared DRAM Systems

Onur Mutlu   Thomas Moscibroda  
Microsoft Research  
{onur,moscitho}@microsoft.com

---

# On PAR-BS

---

- Variants implemented in Samsung SoC memory controllers

Effective platform level approach and DRAM accesses are crucial to system performance. This paper touches this topics and suggest a superior approach to current known techniques.

**Review from ISCA 2008**

---

# ATLAS Memory Scheduler [HPCA'10]

---

- Yoongu Kim, Dongsu Han, Onur Mutlu, and Mor Harchol-Balter, **"ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers"**  
*Proceedings of the 16th International Symposium on High-Performance Computer Architecture (HPCA)*, Bangalore, India, January 2010. [Slides \(pptx\)](#)

## **ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers**

Yoongu Kim   Dongsu Han   Onur Mutlu   Mor Harchol-Balter

Carnegie Mellon University

---



# Thread Cluster Memory Scheduling [MICRO'10]

---

- Yoongu Kim, Michael Papamichael, Onur Mutlu, and Mor Harchol-Balter,

## **"Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior"**

*Proceedings of the 43rd International Symposium on Microarchitecture (**MICRO**), pages 65-76, Atlanta, GA, December 2010. [Slides \(pptx\)](#) [\(pdf\)](#)*

## Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior

Yoongu Kim

yoonguk@ece.cmu.edu

Michael Papamichael

papamix@cs.cmu.edu

Onur Mutlu

onur@cmu.edu

Mor Harchol-Balter

harchol@cs.cmu.edu

Carnegie Mellon University

---

- Lavanya Subramanian, Donghyuk Lee, Vivek Seshadri, Harsha Rastogi, and Onur Mutlu,  
**"The Blacklisting Memory Scheduler: Achieving High Performance and Fairness at Low Cost"**  
*Proceedings of the 32nd IEEE International Conference on Computer Design (ICCD)*, Seoul, South Korea, October 2014.  
[[Slides \(pptx\)](#)] [[pdf](#)]

## The Blacklisting Memory Scheduler: Achieving High Performance and Fairness at Low Cost

Lavanya Subramanian, Donghyuk Lee, Vivek Seshadri, Harsha Rastogi, Onur Mutlu  
Carnegie Mellon University  
{lsubrama,donghyu1,visesh,harshar,onur}@cmu.edu

# Staged Memory Scheduling: CPU-GPU [ISCA'12]

---

- Rachata Ausavarungnirun, Kevin Chang, Lavanya Subramanian, Gabriel Loh, and Onur Mutlu,  
**"Staged Memory Scheduling: Achieving High Performance and Scalability in Heterogeneous Systems"**  
*Proceedings of the 39th International Symposium on Computer Architecture (ISCA)*, Portland, OR, June 2012. [Slides \(pptx\)](#)

## Staged Memory Scheduling: Achieving High Performance and Scalability in Heterogeneous Systems

Rachata Ausavarungnirun<sup>†</sup> Kevin Kai-Wei Chang<sup>†</sup> Lavanya Subramanian<sup>†</sup> Gabriel H. Loh<sup>‡</sup> Onur Mutlu<sup>†</sup>

<sup>†</sup>Carnegie Mellon University  
{rachata,kevincha,lsubrama,onur}@cmu.edu

<sup>‡</sup>Advanced Micro Devices, Inc.  
gabe.loh@amd.com

# DASH: Heterogeneous Systems [TACO'16]

---

- Hiroyuki Usui, Lavanya Subramanian, Kevin Kai-Wei Chang, and Onur Mutlu,  
**"DASH: Deadline-Aware High-Performance Memory Scheduler for Heterogeneous Systems with Hardware Accelerators"**  
*ACM Transactions on Architecture and Code Optimization* (**TACO**), Vol. 12, January 2016.  
Presented at the 11th HiPEAC Conference, Prague, Czech Republic, January 2016.  
[Slides (pptx)] [pdf]  
[Source Code]

## **DASH: Deadline-Aware High-Performance Memory Scheduler for Heterogeneous Systems with Hardware Accelerators**

HIROYUKI USUI, LAVANYA SUBRAMANIAN, KEVIN KAI-WEI CHANG,  
and ONUR MUTLU, Carnegie Mellon University

# MISE: Predictable Performance [HPCA'13]

---

- Lavanya Subramanian, Vivek Seshadri, Yoongu Kim, Ben Jaiyen, and Onur Mutlu,  
**"MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems"**  
*Proceedings of the 19th International Symposium on High-Performance Computer Architecture (**HPCA**), Shenzhen, China, February 2013. Slides (pptx)*

## MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems

Lavanya Subramanian

Vivek Seshadri

Yoongu Kim

Ben Jaiyen

Onur Mutlu

Carnegie Mellon University

# ASM: Predictable Performance [MICRO'15]

---

- Lavanya Subramanian, Vivek Seshadri, Arnab Ghosh, Samira Khan, and Onur Mutlu,  
**"The Application Slowdown Model: Quantifying and Controlling the Impact of Inter-Application Interference at Shared Caches and Main Memory"**  
*Proceedings of the 48th International Symposium on Microarchitecture (MICRO)*, Waikiki, Hawaii, USA, December 2015.  
[[Slides \(pptx\)](#)] [[pdf](#)] [[Lightning Session Slides \(pptx\)](#)] [[pdf](#)] [[Poster \(pptx\)](#)] [[pdf](#)]  
[[Source Code](#)]

## The Application Slowdown Model: Quantifying and Controlling the Impact of Inter-Application Interference at Shared Caches and Main Memory

Lavanya Subramanian\*§      Vivek Seshadri\*      Arnab Ghosh\*†  
Samira Khan\*‡      Onur Mutlu\*

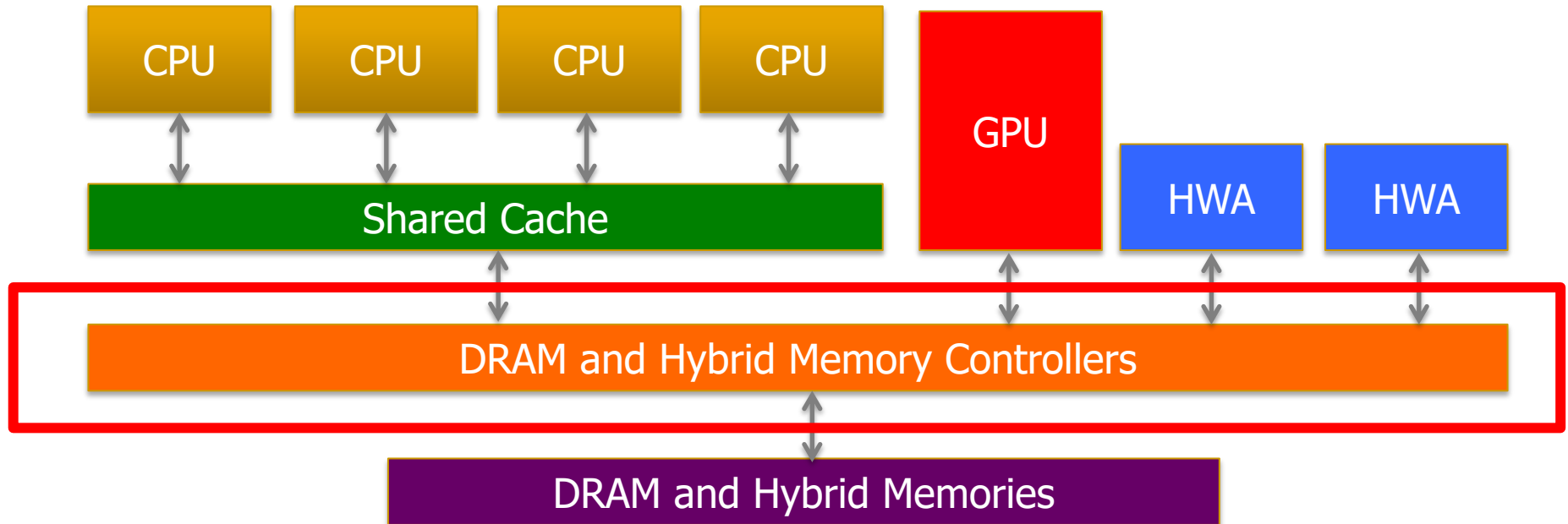
\*Carnegie Mellon University    §Intel Labs    †IIT Kanpur    ‡University of Virginia

Memory Controllers  
are critical to research

They will become  
even more important

# Memory Control is Getting More Complex

---



- Heterogeneous agents: CPUs, GPUs, and HWAs
- Main memory interference between CPUs, GPUs, HWAs

**Many goals, many constraints, many metrics ...**



# Memory Control w/ Machine Learning [ISCA'08]

---

- Engin Ipek, Onur Mutlu, José F. Martínez, and Rich Caruana,  
**"Self Optimizing Memory Controllers: A Reinforcement Learning Approach"**  
*Proceedings of the 35th International Symposium on Computer Architecture (ISCA)*, pages 39-50, Beijing, China, June 2008. [Slides \(pptx\)](#)

## Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek<sup>1,2</sup>   Onur Mutlu<sup>2</sup>   José F. Martínez<sup>1</sup>   Rich Caruana<sup>1</sup>

<sup>1</sup>Cornell University, Ithaca, NY 14850 USA

<sup>2</sup>Microsoft Research, Redmond, WA 98052 USA

## Memory Controllers: Many New Problems

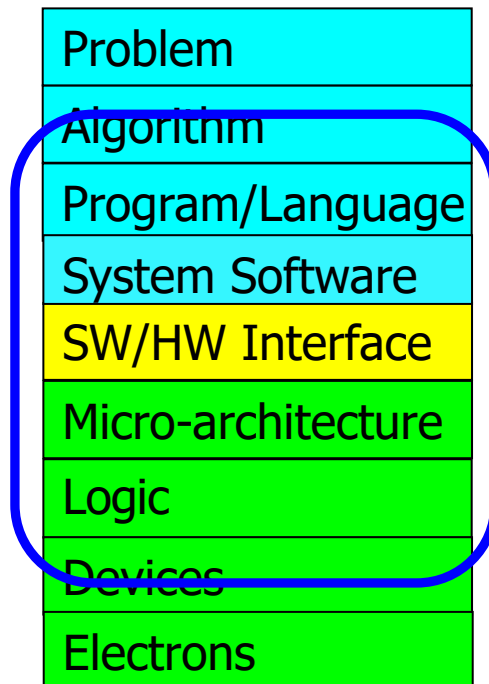
## Main Memory Needs Intelligent Controllers

# We Will See More Examples of This

---

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**

# What We Will Cover In The Next Few Lectures

# Agenda for The Next Few Lectures

---

- Memory Importance and Trends
- RowHammer: Memory Reliability and Security
- Computation in Memory (Processing in/near Memory)
- Low-Latency Memory
- Data-Driven and Data-Aware Architectures
- Memory Controllers and Memory QoS
- Guiding Principles & Research Topics

# An “Early” Position Paper [IMW’13]

---

- Onur Mutlu,  
**"Memory Scaling: A Systems Architecture Perspective"**  
*Proceedings of the 5th International Memory Workshop (**IMW**), Monterey, CA, May 2013. Slides  
(pptx) (pdf)  
EETimes Reprint*

## Memory Scaling: A Systems Architecture Perspective

Onur Mutlu  
Carnegie Mellon University  
onur@cmu.edu  
<http://users.ece.cmu.edu/~omutlu/>

# An Extended Version: Memory Scaling

---

- Onur Mutlu,  
**"Main Memory Scaling: Challenges and Solution Directions"**  
*Invited Book Chapter in More than Moore Technologies for Next Generation Computer Design, pp. 127-153, Springer, 2015.*

## Chapter 6

# Main Memory Scaling: Challenges and Solution Directions

*Onur Mutlu, Carnegie Mellon University*

**Part of your Homework 1 assignment**

---



# Challenges in Memory Scaling

---

- Data retention (need for refresh)
- Reliability and vulnerabilities (e.g., RowHammer)
- Latency and parallelism (e.g., bank conflicts)
- Energy & power
- Memory's inability to do more than store data

# A Recent Retrospective Paper [TCAD'19]

---

- 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

# Computer Architecture

## Lecture 4a: Memory Systems: Solution Directions

Prof. Onur Mutlu

ETH Zürich

Fall 2020

25 September 2020

# Backup Slides

# Readings, Videos, Reference Materials

# Accelerated Memory Course (~6.5 hours)

---

## ■ ACACES 2018

- ❑ Memory Systems and Memory-Centric Computing Systems
- ❑ Taught by Onur Mutlu July 9-13, 2018
- ❑ ~6.5 hours of lectures

## ■ Website for the Course including Videos, Slides, Papers

- ❑ [https://safari.ethz.ch/memory\\_systems/ACACES2018/](https://safari.ethz.ch/memory_systems/ACACES2018/)
- ❑ <https://www.youtube.com/playlist?list=PL5Q2soXY2Zi-HXxomthrpDpMJm05P6J9x>

## ■ All Papers are at:

- ❑ <https://people.inf.ethz.ch/omutlu/projects.htm>
- ❑ Final lecture notes and readings (for all topics)

# Longer Memory Course (~18 hours)

---

## ■ Tu Wien 2019

- ❑ Memory Systems and Memory-Centric Computing Systems
- ❑ Taught by Onur Mutlu June 12-19, 2019
- ❑ ~18 hours of lectures

## ■ Website for the Course including Videos, Slides, Papers

- ❑ [https://safari.ethz.ch/memory\\_systems/TUWien2019](https://safari.ethz.ch/memory_systems/TUWien2019)
- ❑ [https://www.youtube.com/playlist?list=PL5Q2soXY2Zi\\_gntM55VoMIKlw7YrXOhbl](https://www.youtube.com/playlist?list=PL5Q2soXY2Zi_gntM55VoMIKlw7YrXOhbl)

## ■ All Papers are at:

- ❑ <https://people.inf.ethz.ch/omutlu/projects.htm>
- ❑ Final lecture notes and readings (for all topics)

# Related 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)



# Reference Overview Paper I

---

## Processing Data Where It Makes Sense: Enabling In-Memory Computation

Onur Mutlu<sup>a,b</sup>, Saugata Ghose<sup>b</sup>, Juan Gómez-Luna<sup>a</sup>, Rachata Ausavarungnirun<sup>b,c</sup>

<sup>a</sup>*ETH Zürich*

<sup>b</sup>*Carnegie Mellon University*

<sup>c</sup>*King Mongkut's University of Technology North Bangkok*

Onur Mutlu, Saugata Ghose, Juan Gomez-Luna, and Rachata Ausavarungnirun,  
**"Processing Data Where It Makes Sense: Enabling In-Memory  
Computation"**

*Invited paper in Microprocessors and Microsystems (**MICPRO**), June 2019.  
[arXiv version]*

# Reference Overview Paper II

---

## **A Workload and Programming Ease Driven Perspective of Processing-in-Memory**

Saugata Ghose<sup>†</sup>   Amirali Boroumand<sup>†</sup>   Jeremie S. Kim<sup>†§</sup>   Juan Gómez-Luna<sup>§</sup>   Onur Mutlu<sup>§†</sup>

<sup>†</sup>*Carnegie Mellon University*

<sup>§</sup>*ETH Zürich*

Saugata Ghose, Amirali Boroumand, Jeremie S. Kim, Juan Gomez-Luna, and Onur Mutlu,

**"Processing-in-Memory: A Workload-Driven Perspective"**

*Invited Article in IBM Journal of Research & Development, Special Issue on Hardware for Artificial Intelligence, to appear in November 2019.*

[Preliminary arXiv version]

# Reference Overview Paper III

---

## **Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms, Future Research Directions**

SAUGATA GHOSE, KEVIN HSIEH, AMIRALI BOROUMAND,  
RACHATA AUSAVARUNGNIRUN

Carnegie Mellon University

ONUR MUTLU

ETH Zürich and Carnegie Mellon University

Saugata Ghose, Kevin Hsieh, Amirali Boroumand, Rachata Ausavarungnirun, Onur Mutlu,  
**"Enabling the Adoption of Processing-in-Memory: Challenges, Mechanisms,  
Future Research Directions"**

*Invited Book Chapter, to appear in 2018.*

[[Preliminary arxiv.org version](https://arxiv.org/pdf/1802.00320.pdf)]

# Reference Overview Paper IV

---

- Onur Mutlu and Lavanya Subramanian,  
**"Research Problems and Opportunities in Memory Systems"**  
*Invited Article in Supercomputing Frontiers and Innovations (SUPERFRI), 2014/2015.*

## Research Problems and Opportunities in Memory Systems

*Onur Mutlu<sup>1</sup>, Lavanya Subramanian<sup>1</sup>*

# Reference Overview Paper V

---

- Onur Mutlu,  
**"The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser"**  
*Invited Paper in Proceedings of the Design, Automation, and Test in Europe Conference (**DATE**), Lausanne, Switzerland, March 2017.*  
[[Slides \(pptx\)](#) ([pdf](#))]

## The RowHammer Problem and Other Issues We May Face as Memory Becomes Denser

Onur Mutlu  
ETH Zürich  
[onur.mutlu@inf.ethz.ch](mailto:onur.mutlu@inf.ethz.ch)  
<https://people.inf.ethz.ch/omutlu>

# Reference Overview Paper VI

---

- Onur Mutlu,  
**"Memory Scaling: A Systems Architecture Perspective"**

*Technical talk at MemCon 2013 (**MEMCON**), Santa Clara, CA, August 2013. [[Slides \(pptx\)](#)] [[pdf](#)]  
[[Video](#)] [[Coverage on StorageSearch](#)]*

## Memory Scaling: A Systems Architecture Perspective

Onur Mutlu  
Carnegie Mellon University  
onur@cmu.edu  
<http://users.ece.cmu.edu/~omutlu/>



*Proceedings of the IEEE, Sept. 2017*

## Error Characterization, Mitigation, and Recovery in Flash-Memory-Based Solid-State Drives

*This paper reviews the most recent advances in solid-state drive (SSD) error characterization, mitigation, and data recovery techniques to improve both SSD's reliability and lifetime.*

By YU CAI, SAUGATA GHOSE, ERICH F. HARATSCH, YIXIN LUO, AND ONUR MUTLU

# Reference Overview Paper VIII

---

- 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



# Reference Overview Paper IX

---

- Vivek Seshadri and Onur Mutlu,  
**"In-DRAM Bulk Bitwise Execution Engine"**  
*Invited Book Chapter in Advances in Computers*, to appear  
in 2020.  
[[Preliminary arXiv version](#)]

## In-DRAM Bulk Bitwise Execution Engine

Vivek Seshadri  
Microsoft Research India  
`visesha@microsoft.com`

Onur Mutlu  
ETH Zürich  
`onur.mutlu@inf.ethz.ch`

# Related Videos and Course Materials (I)

---

- **Undergraduate Digital Design & Computer Architecture Course Lecture Videos (2020, 2019, 2018, 2017, 2015, 2014, 2013)**
- **Undergraduate Digital Design & Computer Architecture Course Materials (2020, 2019, 2018, 2015, 2014, 2013)**
- **Graduate Computer Architecture Course Lecture Videos (2019, 2018, 2017, 2015, 2013)**
- **Graduate Computer Architecture Course Materials (2019, 2018, 2017, 2015, 2013)**
- **Parallel Computer Architecture Course Materials (Lecture Videos)**

# Related Videos and Course Materials (II)

---

- **Seminar in Computer Architecture Course Lecture Videos (Spring 2020, Fall 2019, Spring 2019, 2018)**
- **Seminar in Computer Architecture Course Materials (Spring 2020, Fall 2019, Spring 2019, 2018)**
- **Memory Systems Course Lecture Videos (Sept 2019, July 2019, June 2019, October 2018)**
- **Memory Systems Short Course Lecture Materials (Sept 2019, July 2019, June 2019, October 2018)**
- **ACACES Summer School Memory Systems Course Lecture Videos (2018, 2013)**
- **ACACES Summer School Memory Systems Course Materials (2018, 2013)**

# Some Open Source Tools (I)

---

- Rowhammer – Program to Induce RowHammer Errors
  - <https://github.com/CMU-SAFARI/rowhammer>
- Ramulator – Fast and Extensible DRAM Simulator
  - <https://github.com/CMU-SAFARI/ramulator>
- MemSim – Simple Memory Simulator
  - <https://github.com/CMU-SAFARI/memsim>
- NOCulator – Flexible Network-on-Chip Simulator
  - <https://github.com/CMU-SAFARI/NOCulator>
- SoftMC – FPGA-Based DRAM Testing Infrastructure
  - <https://github.com/CMU-SAFARI/SoftMC>
- Other open-source software from my group
  - <https://github.com/CMU-SAFARI/>
  - <http://www.ece.cmu.edu/~safari/tools.html>

# Some Open Source Tools (II)

---

- MQSim – A Fast Modern SSD Simulator
  - <https://github.com/CMU-SAFARI/MQSim>
- Mosaic – GPU Simulator Supporting Concurrent Applications
  - <https://github.com/CMU-SAFARI/Mosaic>
- IMPICA – Processing in 3D-Stacked Memory Simulator
  - <https://github.com/CMU-SAFARI/IMPICA>
- SMLA – Detailed 3D-Stacked Memory Simulator
  - <https://github.com/CMU-SAFARI/SMLA>
- HWASim – Simulator for Heterogeneous CPU-HWA Systems
  - <https://github.com/CMU-SAFARI/HWASim>
- Other open-source software from my group
  - <https://github.com/CMU-SAFARI/>
  - <http://www.ece.cmu.edu/~safari/tools.html>

# More Open Source Tools (III)

- A lot more open-source software from my group
  - ❑ <https://github.com/CMU-SAFARI/>
  - ❑ <http://www.ece.cmu.edu/~safari/tools.html>

The screenshot shows the GitHub repository page for the SAFARI Research Group. The header includes the SAFARI logo and the text "SAFARI Research Group at ETH Zurich and Carnegie Mellon University". Below this, it states "Site for source code and tools distribution from SAFARI Research Group at ETH Zurich and Carnegie Mellon University." and provides contact information: "ETH Zurich and Carnegi...", "http://www.ece.cmu.ed...", and "omutlu@gmail.com". The repository statistics bar shows 30 Repositories, 27 People, 1 Teams, and 0 Projects. The main content area features a search bar, filters for Type and Language, and a "New" button. The repository "MQSim" is highlighted, with a description: "MQSim is a fast and accurate simulator modeling the performance of modern multi-queue (MQ) SSDs as well as traditional SATA based SSDs. MQSim faithfully models new high-bandwidth protocol implementations, steady-state SSD conditions, and the full end-to-end latency of requests in modern SSDs. It is described in detail in the FAST 2018 paper by A...". The repository has 14 stars, 14 forks, and is maintained by MIT. The "Top languages" section lists C++, C, C#, AGS Script, and Verilog. The "Most used topics" section lists dram and reliability.

**SAFARI** SAFARI Research Group at ETH Zurich and Carnegie Mellon University

Site for source code and tools distribution from SAFARI Research Group at ETH Zurich and Carnegie Mellon University.

ETH Zurich and Carnegi... http://www.ece.cmu.ed... omutlu@gmail.com

Repositories 30 People 27 Teams 1 Projects 0 Settings

Search repositories... Type: All Language: All Customize pinned repositories New

**MQSim**

MQSim is a fast and accurate simulator modeling the performance of modern multi-queue (MQ) SSDs as well as traditional SATA based SSDs. MQSim faithfully models new high-bandwidth protocol implementations, steady-state SSD conditions, and the full end-to-end latency of requests in modern SSDs. It is described in detail in the FAST 2018 paper by A...

C++ ★ 14 14 MIT Updated 8 days ago

Top languages

- C++
- C
- C#
- AGS Script
- Verilog

Most used topics

- dram
- reliability

Manage

## ramulator-pim

A fast and flexible simulation infrastructure for exploring general-purpose processing-in-memory (PIM) architectures. Ramulator-PIM combines a widely-used simulator for out-of-order and in-order processors (ZSim) with Ramulator, a DRAM simulator with memory models for DDRx, LPDDRx, GDDRx, WIOx, HBMx, and HMCx. Ramulator is described in the IEEE ...

● C++ 🍴 11 ☆ 29 ⓘ 6 📄 0 Updated 19 days ago

## SMASH

SMASH is a hardware-software cooperative mechanism that enables highly-efficient indexing and storage of sparse matrices. The key idea of SMASH is to compress sparse matrices with a hierarchical bitmap compression format that can be accelerated from hardware.

Described by Kanellopoulos et al. (MICRO '19)  
<https://people.inf.ethz.ch/omutlu/pub/SMA...>

● C 🍴 1 ☆ 6 ⓘ 0 📄 0 Updated on May 17

## MQSim

MQSim is a fast and accurate simulator modeling the performance of modern multi-queue (MQ) SSDs as well as traditional SATA based SSDs. MQSim faithfully models new high-bandwidth protocol implementations, steady-state SSD conditions, and the full end-to-end latency of requests in modern SSDs. It is described in detail in the FAST 2018 paper by A...

● C++ 🍴 MIT 🍴 54 ☆ 62 ⓘ 10 📄 1 Updated on May 15

## Apollo

Apollo is an assembly polishing algorithm that attempts to correct the errors in an assembly. It can take multiple set of reads in a single run and polish the assemblies of genomes of any size. Described in the Bioinformatics journal paper (2020) by Firtina et al. at  
<https://people.inf.ethz.ch/omutlu/pub/apollo-technology-independent-genome-assem...>

● C++ 🍴 GPL-3.0 🍴 1 ☆ 12 ⓘ 0 📄 0 Updated on May 10

## ramulator

A Fast and Extensible DRAM Simulator, with built-in support for modeling many different DRAM technologies including DDRx, LPDDRx, GDDRx, WIOx, HBMx, and various academic proposals. Described in the IEEE CAL 2015 paper by Kim et al. at  
[http://users.ece.cmu.edu/~omutlu/pub/ramulator\\_dram\\_simulator-ieee-cal15.pdf](http://users.ece.cmu.edu/~omutlu/pub/ramulator_dram_simulator-ieee-cal15.pdf)

● C++ 🍴 MIT 🍴 93 ☆ 170 ⓘ 37 📄 2 Updated on Apr 13

## Shifted-Hamming-Distance

Source code for the Shifted Hamming Distance (SHD) filtering mechanism for sequence alignment. Described in the Bioinformatics journal paper (2015) by Xin et al. at  
[http://users.ece.cmu.edu/~omutlu/pub/shifted-hamming-distance\\_bioinformatics15\\_proofs.pdf](http://users.ece.cmu.edu/~omutlu/pub/shifted-hamming-distance_bioinformatics15_proofs.pdf)

● C 🍴 GPL-2.0 🍴 5 ☆ 20 ⓘ 0 📄 1 Updated on Mar 29

## SneakySnake

The first and the only pre-alignment filtering algorithm that works on all modern high-performance computing architectures. It works efficiently and fast on CPU, FPGA, and GPU architectures and that greatly (by more than two orders of magnitude) expedites sequence alignment calculation. Described by Alser et al. (preliminary version at <https://a...>)

● VHDL 🍴 GPL-3.0 🍴 3 ☆ 11 ⓘ 0 📄 0 Updated on Mar 10

## AirLift

AirLift is a tool that updates mapped reads from one reference genome to another. Unlike existing tools, It accounts for regions not shared between the two reference genomes and enables remapping across all parts of the references. Described by Kim et al. (preliminary version at <http://arxiv.org/abs/1912.08735>)

● C 🍴 0 ☆ 3 ⓘ 0 📄 0 Updated on Feb 19

## GPGPUSim-Ramulator

The source code for GPGPUSim+Ramulator simulator. In this version, GPGPUSim uses Ramulator to simulate the DRAM. This simulator is used to produce some of the

# Referenced Papers and Talks

---

- All are available at

**<https://people.inf.ethz.ch/omutlu/projects.htm>**

**<http://scholar.google.com/citations?user=7XyGUGkAAAAJ&hl=en>**

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



# 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\\_5MGV6EnXEJHnV2YFBJl&index=15](https://www.youtube.com/watch?v=tcQ3zZ3JpuA&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=15)

# Ramulator: A Fast and Extensible DRAM Simulator

**[IEEE Comp Arch Letters'15]**

# Ramulator Motivation

- DRAM and Memory Controller landscape is changing
- Many new and upcoming standards
- Many new controller designs
- A fast and easy-to-extend simulator is very much needed

<i>Segment</i>	<i>DRAM Standards &amp; Architectures</i>
Commodity	DDR3 (2007) [14]; DDR4 (2012) [18]
Low-Power	LPDDR3 (2012) [17]; LPDDR4 (2014) [20]
Graphics	GDDR5 (2009) [15]
Performance	eDRAM [28], [32]; RLDram3 (2011) [29]
3D-Stacked	WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (2015) [13]; HBM (2013) [19]; HMC1.0 (2013) [10]; HMC1.1 (2014) [11]
Academic	SBA/SSA (2010) [38]; Staged Reads (2012) [8]; RAIDR (2012) [27]; SALP (2012) [24]; TL-DRAM (2013) [26]; RowClone (2013) [37]; Half-DRAM (2014) [39]; Row-Buffer Decoupling (2014) [33]; SARP (2014) [6]; AL-DRAM (2015) [25]

Table 1. Landscape of DRAM-based memory

# Ramulator

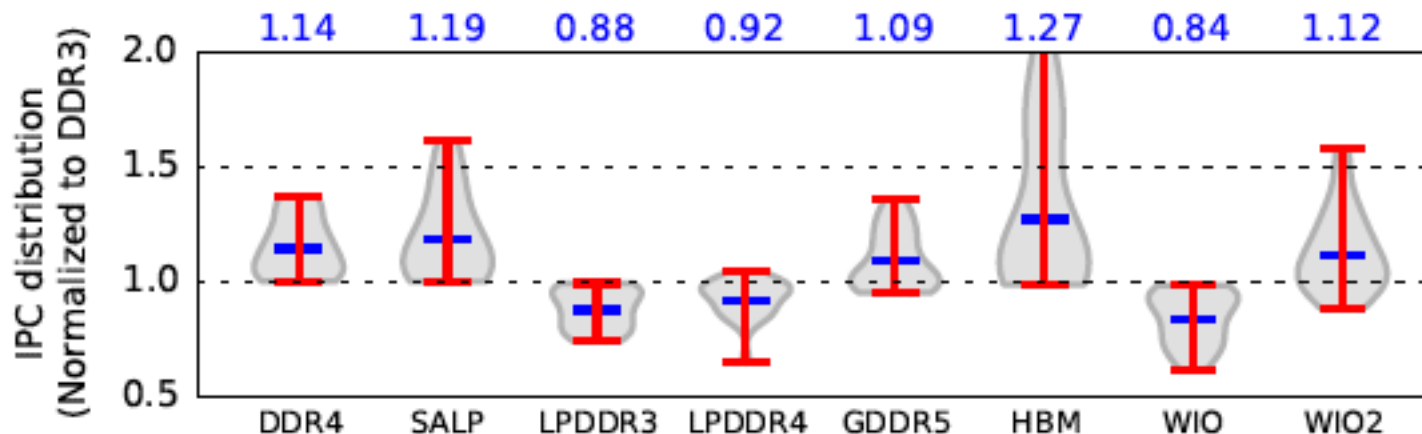
- Provides out-of-the box support for many DRAM standards:
  - DDR3/4, LPDDR3/4, GDDR5, WIO1/2, HBM, plus new proposals (SALP, AL-DRAM, TLDRAM, RowClone, and SARP)
- ~2.5X faster than fastest open-source simulator
- Modular and extensible to different standards

<i>Simulator</i> (clang -O3)	<i>Cycles (10<sup>6</sup>)</i>		<i>Runtime (sec.)</i>		<i>Req/sec (10<sup>3</sup>)</i>		<i>Memory</i> (MB)
	<i>Random</i>	<i>Stream</i>	<i>Random</i>	<i>Stream</i>	<i>Random</i>	<i>Stream</i>	
Ramulator	652	411	752	249	133	402	2.1
DRAMSim2	645	413	2,030	876	49	114	1.2
USIMM	661	409	1,880	750	53	133	4.5
DrSim	647	406	18,109	12,984	6	8	1.6
NVMain	666	413	6,881	5,023	15	20	4,230.0

Table 3. Comparison of five simulators using two traces

# Case Study: Comparison of DRAM Standards

<i>Standard</i>	<i>Rate (MT/s)</i>	<i>Timing (CL-RCD-RP)</i>	<i>Data-Bus (Width×Chan.)</i>	<i>Rank-per-Chan</i>	<i>BW (GB/s)</i>
DDR3	1,600	11-11-11	64-bit × 1	1	11.9
DDR4	2,400	16-16-16	64-bit × 1	1	17.9
SALP <sup>†</sup>	1,600	11-11-11	64-bit × 1	1	11.9
LPDDR3	1,600	12-15-15	64-bit × 1	1	11.9
LPDDR4	2,400	22-22-22	32-bit × 2*	1	17.9
GDDR5 [12]	6,000	18-18-18	64-bit × 1	1	44.7
HBM	1,000	7-7-7	128-bit × 8*	1	119.2
WIO	266	7-7-7	128-bit × 4*	1	15.9
WIO2	1,066	9-10-10	128-bit × 8*	1	127.2



Across 22 workloads, simple CPU model

Figure 2. Performance comparison of DRAM standards

# Ramulator Paper and Source Code

---

- Yoongu Kim, Weikun Yang, and Onur Mutlu,  
**"Ramulator: A Fast and Extensible DRAM Simulator"**  
*IEEE Computer Architecture Letters (CAL)*, March 2015.  
[[Source Code](#)]
- Source code is released under the liberal MIT License
  - <https://github.com/CMU-SAFARI/ramulator>

## Ramulator: A Fast and Extensible DRAM Simulator

Yoongu Kim<sup>1</sup>      Weikun Yang<sup>1,2</sup>      Onur Mutlu<sup>1</sup>  
<sup>1</sup>Carnegie Mellon University      <sup>2</sup>Peking University

# Optional Assignment

---

- Review the Ramulator paper
  - Email me your review ([omutlu@gmail.com](mailto:omutlu@gmail.com))
- Download and run Ramulator
  - Compare DDR3, DDR4, SALP, HBM for the libquantum benchmark (provided in Ramulator repository)
  - Email me your report ([omutlu@gmail.com](mailto:omutlu@gmail.com))
- This **will** help you get into **memory systems research**

# End of Backup Slides