

Computer Architecture

Lecture 12: Memory Controllers

Ataberk Olgun

Prof. Onur Mutlu

ETH Zürich

Fall 2023

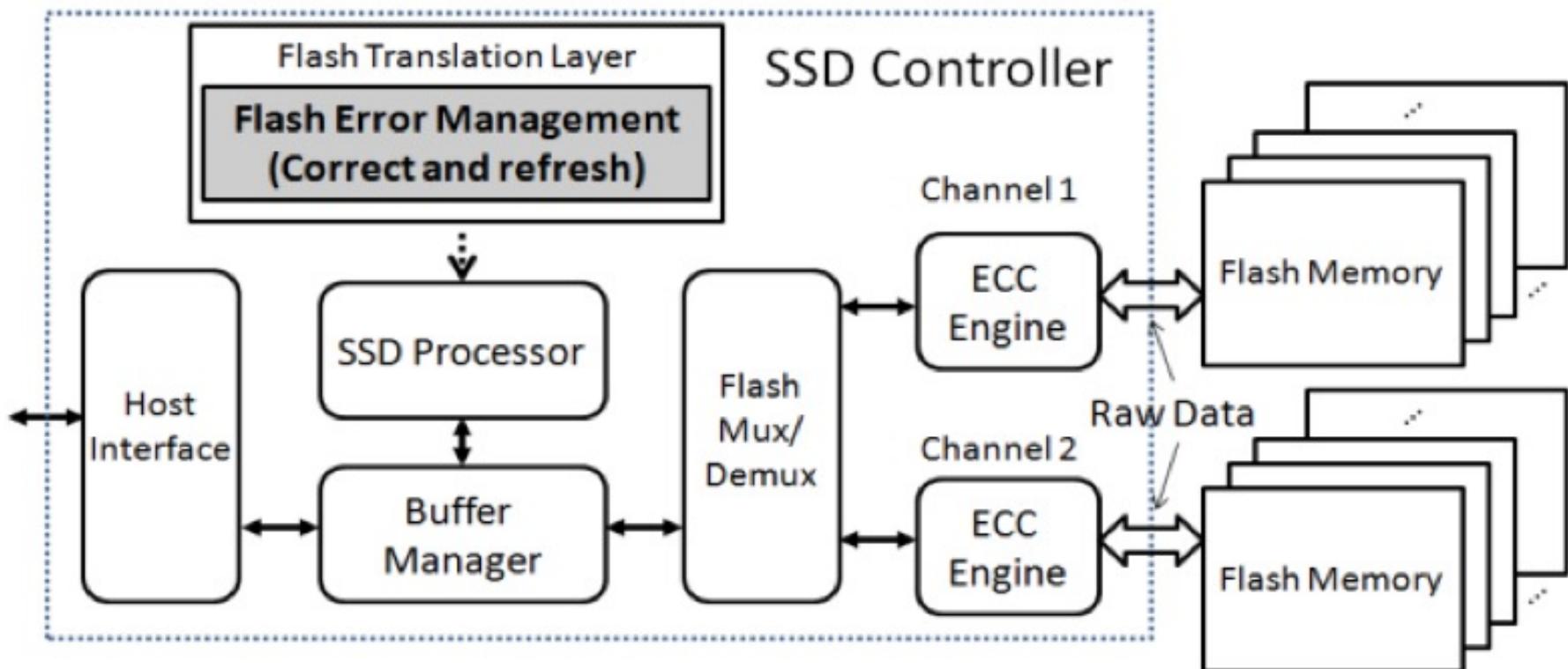
3 November 2023

DRAM versus Other Types of Memories

- Long latency memories have similar characteristics that need to be controlled.
- This lecture will use DRAM as an example, but many scheduling and control issues are similar in the design of controllers for other types of memories
 - Flash memory
 - Other emerging memory technologies
 - Phase Change Memory
 - Spin-Transfer Torque Magnetic Memory
 - These other technologies can also place other demands on the controller

Flash Memory (SSD) Controllers

- Similar to DRAM memory controllers, except:
 - They are flash memory specific
 - They do much more: complex error correction, wear leveling, voltage optimization, garbage collection, page remapping, ...



Another View of the SSD Controller

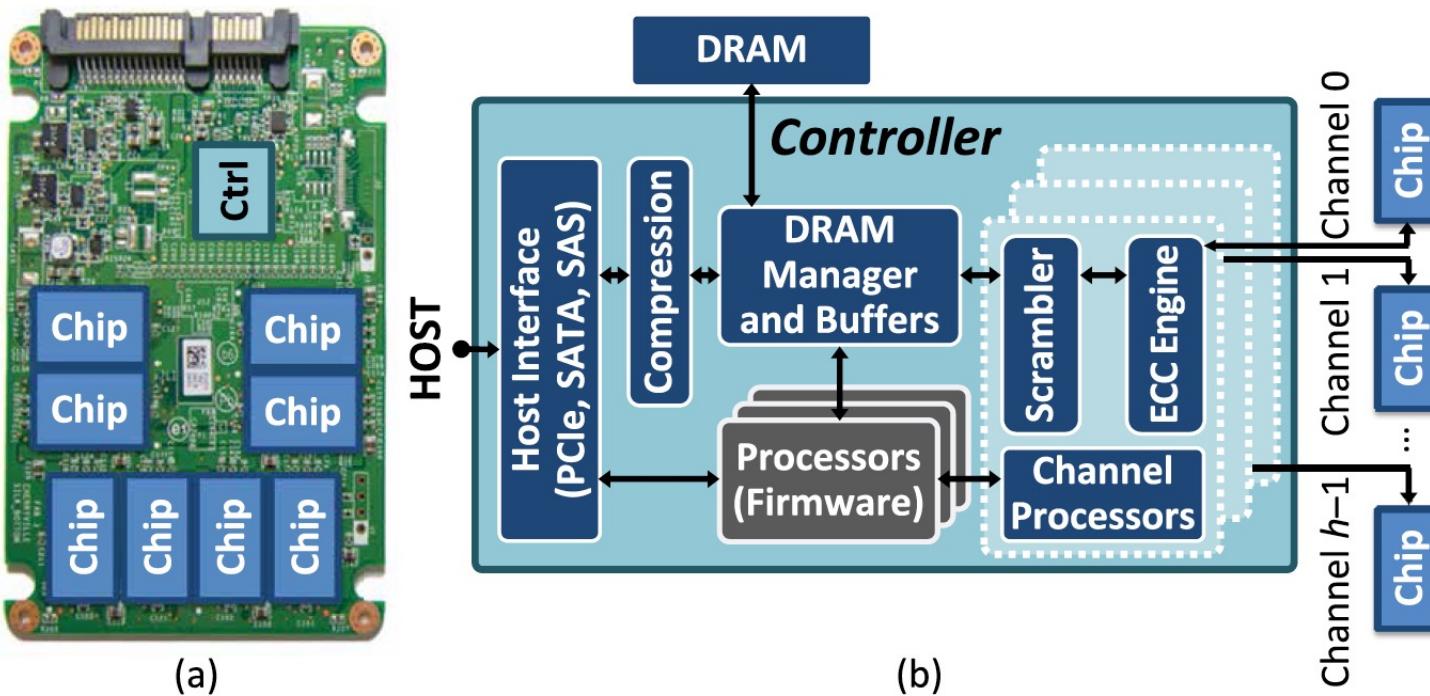


Fig. 1. (a) SSD system architecture, showing controller (Ctrl) and chips. (b) Detailed view of connections between controller components and chips.

On Modern SSD Controllers (I)



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

Many Errors and Their Mitigation [PIEEE'17]

Table 3 List of Different Types of Errors Mitigated by NAND Flash Error Mitigation Mechanisms

Mitigation Mechanism	Error Type				
	P/E Cycling [32,33,42] (§IV-A)	Program [40,42,53] (§IV-B)	Cell-to-Cell Interference [32,35,36,55] (§IV-C)	Data Retention [20,32,34,37,39] (§IV-D)	Read Disturb [20,32,38,62] (§IV-E)
Shadow Program Sequencing [35,40] (Section V-A)		X			
Neighbor-Cell Assisted Error Correction [36] (Section V-B)		X			
Refresh [34,39,67,68] (Section V-C)				X	X
Read-Retry [33,72,107] (Section V-D)	X			X	X
Voltage Optimization [37,38,74] (Section V-E)	X			X	X
Hot Data Management [41,63,70] (Section V-F)	X	X	X	X	X
Adaptive Error Mitigation [43,65,77,78,82] (Section V-G)	X	X	X	X	X

Cai+, "Error Characterization, Mitigation, and Recovery in Flash Memory Based Solid State Drives," Proc. IEEE 2017.

More Up-to-date Version

- Yu Cai, Saugata Ghose, Erich F. Haratsch, Yixin Luo, and Onur Mutlu,
"Errors in Flash-Memory-Based Solid-State Drives: Analysis, Mitigation, and Recovery"

Invited Book Chapter in Inside Solid State Drives, 2018.

[Preliminary arxiv.org version]

Errors in Flash-Memory-Based Solid-State Drives: Analysis, Mitigation, and Recovery

YU CAI, SAUGATA GHOSE

Carnegie Mellon University

ERICH F. HARATSCH

Seagate Technology

YIXIN LUO

Carnegie Mellon University

ONUR MUTLU

ETH Zürich and Carnegie Mellon University

On Modern SSD Controllers (II)

- Arash Tavakkol, Juan Gomez-Luna, Mohammad Sadrosadati, Saugata Ghose, and Onur Mutlu,
"MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices"

Proceedings of the 16th USENIX Conference on File and Storage Technologies (FAST), Oakland, CA, USA, February 2018.

[Slides (pptx) (pdf)]

[Source Code]

MQSim: A Framework for Enabling Realistic Studies of Modern Multi-Queue SSD Devices

Arash Tavakkol[†], Juan Gómez-Luna[†], Mohammad Sadrosadati[†], Saugata Ghose[‡], Onur Mutlu^{†‡}
[†]*ETH Zürich* [‡]*Carnegie Mellon University*

On Modern SSD Controllers (III)

- Arash Tavakkol, Mohammad Sadrosadati, Saugata Ghose, Jeremie Kim, Yixin Luo, Yaohua Wang, Nika Mansouri Ghiasi, Lois Orosa, Juan G. Luna and Onur Mutlu,

"FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives"

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]

FLIN: Enabling Fairness and Enhancing Performance in Modern NVMe Solid State Drives

Arash Tavakkol[†] Mohammad Sadrosadati[†] Saugata Ghose[‡] Jeremie S. Kim^{‡†} Yixin Luo[‡]
Yaohua Wang^{†§} Nika Mansouri Ghiasi[†] Lois Orosa^{†*} Juan Gómez-Luna[†] Onur Mutlu^{†‡}
[†]*ETH Zürich* [‡]*Carnegie Mellon University* [§]*NUDT* ^{*}*Unicamp*

On Modern SSD Controllers (IV)

- Myungsuk Kim, Jisung Park, Geonhee Cho, Yoona Kim, Lois Orosa, Onur Mutlu, and Jihong Kim,

"Evanesco: Architectural Support for Efficient Data Sanitization in Modern Flash-Based Storage Systems"

Proceedings of the 25th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Lausanne, Switzerland, March 2020.

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

[[Talk Video](#) (20 mins)]

Evanesco: Architectural Support for Efficient Data Sanitization in Modern Flash-Based Storage Systems

Myungsuk Kim*

morssola75@davinci.snu.ac.kr
Seoul National University

Jisung Park*

jisung.park@inf.ethz.ch
ETH Zürich & Seoul
National University

Geonhee Cho

ghcho@davinci.snu.ac.kr
Seoul National University

Yoona Kim

yoonakim@davinci.snu.ac.kr
Seoul National University

Lois Orosa
lois.orosa@inf.ethz.ch
ETH Zürich

Onur Mutlu
omutlu@gmail.com
ETH Zürich

Jihong Kim
jihong@davinci.snu.ac.kr
Seoul National University

On Modern SSD Controllers (V)

- Jisung Park, Myungsuk Kim, Myoungjun Chun, Lois Orosa, Jihong Kim, and Onur Mutlu,

"Reducing Solid-State Drive Read Latency by Optimizing Read-Retry"

Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), Virtual, March-April 2021.

[[2-page Extended Abstract](#)]

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

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

[[Short Talk Video](#) (5 mins)]

[[Full Talk Video](#) (19 mins)]

Reducing Solid-State Drive Read Latency by Optimizing Read-Retry

Jisung Park¹ Myungsuk Kim^{2,3} Myoungjun Chun² Lois Orosa¹ Jihong Kim² Onur Mutlu¹

¹ETH Zürich
Switzerland

²Seoul National University
Republic of Korea

³Kyungpook National University
Republic of Korea

On Modern SSD Controllers (VI)

- Yixin Luo, Saugata Ghose, Yu Cai, Erich F. Haratsch, and Onur Mutlu,
"HeatWatch: Improving 3D NAND Flash Memory Device Reliability by Exploiting Self-Recovery and Temperature-Awareness"
Proceedings of the 24th International Symposium on High-Performance Computer Architecture (HPCA), Vienna, Austria, February 2018.
[[Lightning Talk Video](#)] [[Slides \(pptx\)](#) ([pdf](#))] [[Lightning Session Slides \(pptx\)](#) ([pdf](#))]

HeatWatch: Improving 3D NAND Flash Memory Device Reliability by Exploiting Self-Recovery and Temperature Awareness

Yixin Luo[†] Saugata Ghose[†] Yu Cai[‡] Erich F. Haratsch[‡] Onur Mutlu^{§†}

[†]*Carnegie Mellon University*

[‡]*Seagate Technology*

[§]*ETH Zürich*

On Modern SSD Controllers (VII)

- Rakesh Nadig, Mohammad Sadrosadati, Haiyu Mao, Nika Mansouri Ghiasi, Arash Tavakkol, Jisung Park, Hamid Sarbazi-Azad, Juan Gómez Luna, and Onur Mutlu,

"Venice: Improving Solid-State Drive Parallelism at Low Cost via Conflict-Free Accesses"

Proceedings of the 50th International Symposium on Computer Architecture (ISCA), Orlando, FL, USA, June 2023.

[[arXiv version](#)]

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

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

[[Lightning Talk Video](#) (3 minutes)]

[[Talk Video](#) (14 minutes, including Q&A)]

Venice: Improving Solid-State Drive Parallelism at Low Cost via Conflict-Free Accesses

*Rakesh Nadig[§] *Mohammad Sadrosadati[§] Haiyu Mao[§] Nika Mansouri Ghiasi[§]
Arash Tavakkol[§] Jisung Park[§] [▼] Hamid Sarbazi-Azad^{†‡} Juan Gómez Luna[§] Onur Mutlu[§]

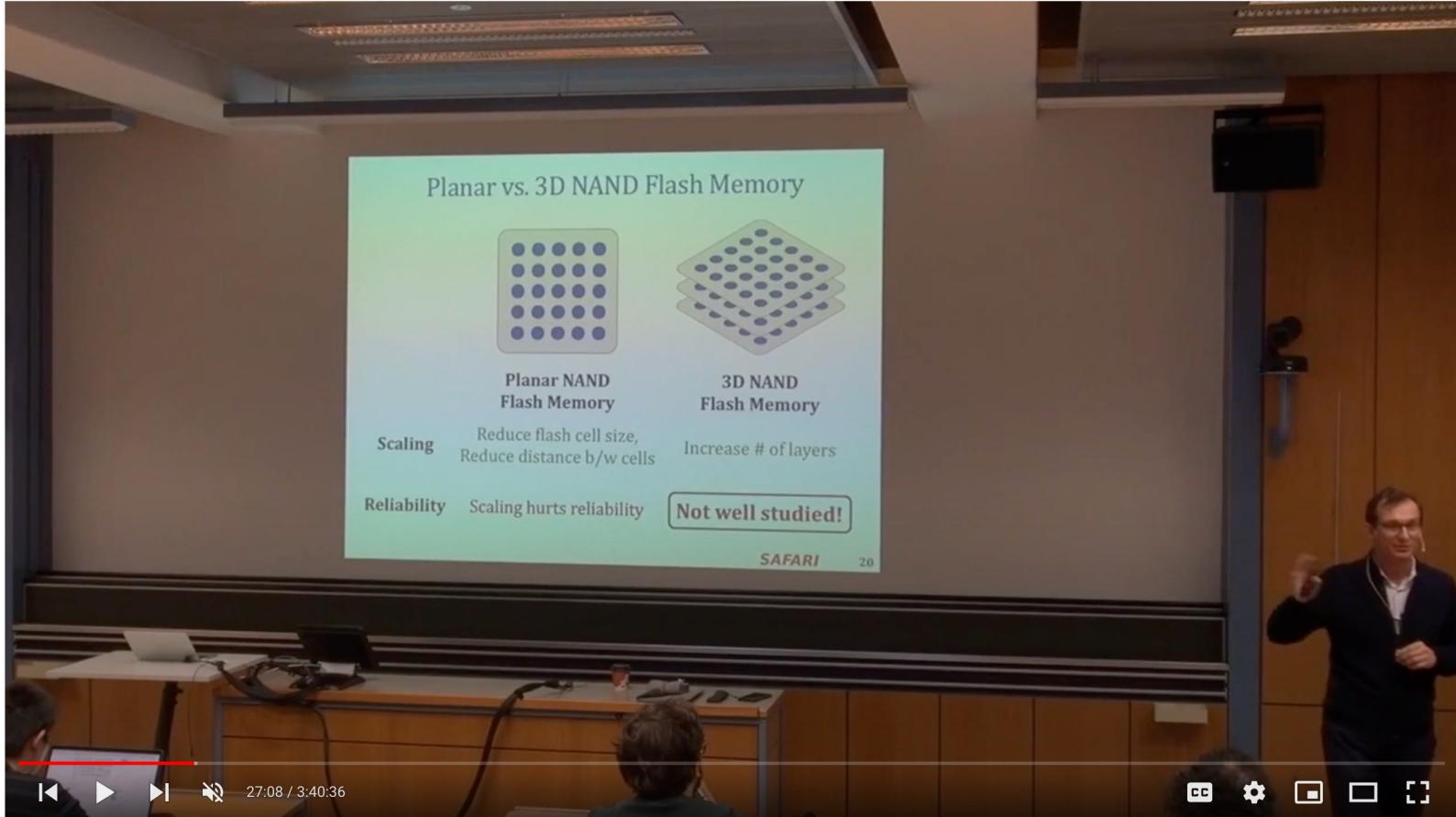
[§]*ETH Zürich*

[▼]*POSTECH*

[†]*Sharif University of Technology*

[‡]*IPM*

Lecture on Flash Memory & SSDs



ETH ZÜRICH HAUPTGEBÄUDE

Computer Architecture - Lecture 26: Flash Memory and Solid-State Drives (ETH Zürich, Fall 2020)

1,771 views • Dec 31, 2020

43 0 SHARE SAVE ...



Onur Mutlu Lectures
19.7K subscribers

ANALYTICS

EDIT VIDEO

Special Course on Flash Memory & SSDs

The image shows a YouTube video player interface. At the top right, there is a small video thumbnail of a man wearing headphones, with the name "Seyyedmoha..." below it. The main video frame contains the title "P&S Modern SSDs" in red, followed by "Basics of NAND Flash-Based SSDs". Below the title, the names "Dr. Mohammad Sadrosadati" and "Prof. Onur Mutlu" are listed, along with "ETH Zürich" and "Spring 2023". The date "17 March 2023" is also mentioned. The video player has a progress bar at the bottom left showing "1:48 / 46:56". To the right of the video frame are various control icons for sharing, saving, and zooming. The overall background is white with a thin yellow border around the video frame.

Modern Solid-State Drives (SSDs) - Lecture 1: Basics of NAND Flash-Based SSDs (Spring 2023)

Onur Mutlu Lectures 36.6K subscribers

Subscribe

1.9K views Streamed 7 months ago Livestream - P&S Modern SSDs (Spring 2023)

Project and Seminars Course: Understanding and Designing Modern NAND Flash-Based SSDs (Solid-State Drives), ETH Zürich, Spring 2023
(https://safari.ethz.ch/projects_and_s...) ...more

Solid-State Drives Course (Spring 2023)

■ Spring 2023 Edition:

- https://safari.ethz.ch/projects_and_seminars/spring2023/doku.php?id=modern_ssds

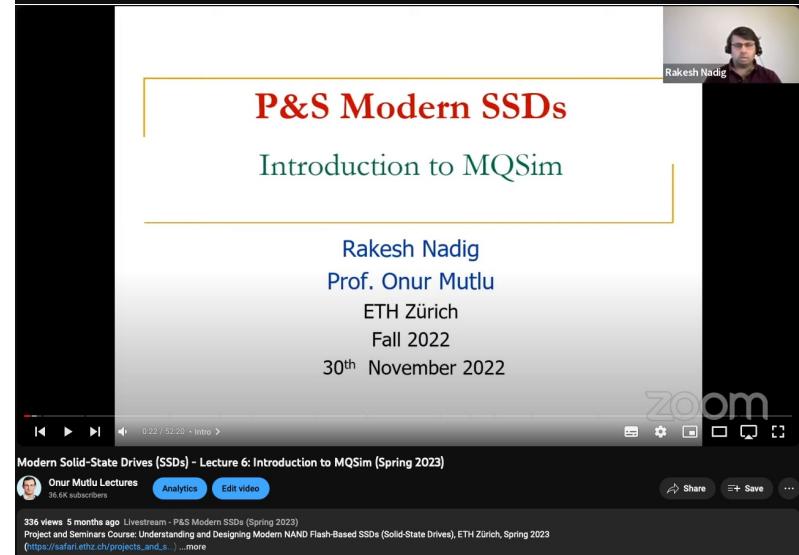
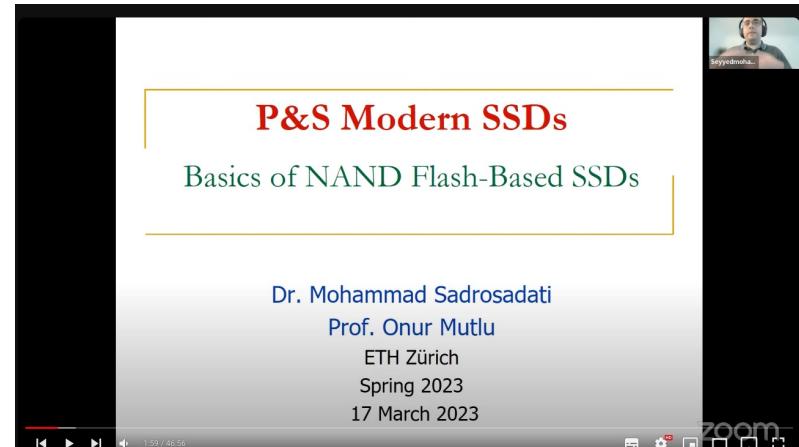
■ Youtube Livestream:

- https://www.youtube.com/watch?v=4VTwOMmsnJY&list=PL5Q2soXY2Zi_8qOM5Icpp8hB2SHtm4z57

■ Project course

- Taken by Bachelor's/Master's students
- SSD Basics and Advanced Topics
- Hands-on research exploration
- Many research readings

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



DRAM Types

- DRAM has different types with different interfaces optimized for different purposes
 - Commodity: DDR, DDR2, DDR3, DDR4, DDR5, ...
 - Low power (for mobile): LPDDR1, ..., LPDDR5, ...
 - High bandwidth (for graphics): GDDR2, ..., GDDR5, ...
 - Low latency: eDRAM, RLDRAM, ...
 - 3D stacked: WIO, HBM, HMC, HBM2.0, ...
 - ...
- Underlying microarchitecture is fundamentally the same
- A flexible memory controller can support various DRAM types
- This complicates the memory controller
 - Difficult to support all types (and upgrades)
 - Analog interface is different for different DRAM types

DRAM Types (circa 2015)

<i>Segment</i>	<i>DRAM Standards & 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

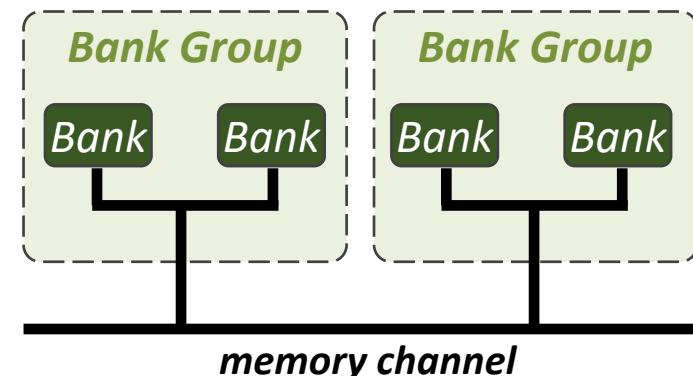
Kim+, "Ramulator: A Flexible and Extensible DRAM Simulator", IEEE CAL 2015.

Modern DRAM Types: Comparison to DDR3

SAFARI

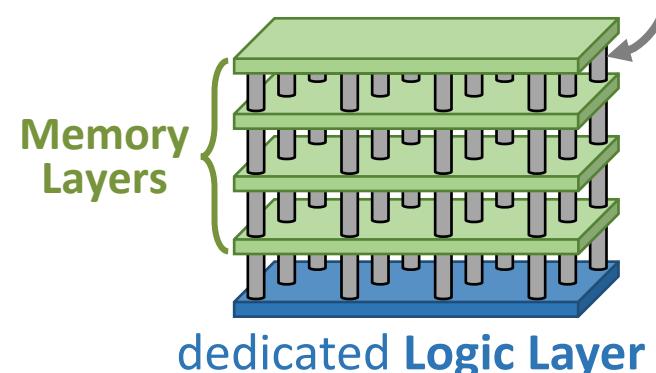
DRAM Type	Banks per Rank	Bank Groups	3D-Stacked	Low-Power
DDR3	8			
DDR4	16	✓	<i>increased latency</i>	
GDDR5	16	✓	<i>increased area/power</i>	
HBM High-Bandwidth Memory	16		✓	
HMC Hybrid Memory Cube	256	<i>narrower rows, higher latency</i>		✓
Wide I/O	4		✓	✓
Wide I/O 2	8		✓	✓
LPDDR3	8			✓
LPDDR4	16			✓

■ Bank groups



■ 3D-stacked DRAM

high bandwidth with Through-Silicon Vias (TSVs)



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 Kim¹ Weikun Yang^{1,2} Onur Mutlu¹
¹Carnegie Mellon University ²Peking University

Ramulator 2.0

- Haocong Luo, Yahya Can Tugrul, F. Nisa Bostancı, Ataberk Olgun, A. Giray Yağlıkçı, and Onur Mutlu,
"Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator"

Preprint on arxiv, August 2023.

[[arXiv version](#)]

[[Ramulator 2.0 Source Code](#)]

- Source code is released under the liberal MIT License
 - <https://github.com/CMU-SAFARI/ramulator2>

Ramulator 2.0: A Modern, Modular, and Extensible DRAM Simulator

Haocong Luo, Yahya Can Tuğrul, F. Nisa Bostancı, Ataberk Olgun, A. Giray Yağlıkçı, and Onur Mutlu

DRAM Types vs. Workloads

- Saugata Ghose, Tianshi Li, Nastaran Hajinazar, Damla Senol Cali, and Onur Mutlu,
"Demystifying Workload–DRAM Interactions: An Experimental Study"
*Proceedings of the ACM International Conference on Measurement and Modeling
of Computer Systems (SIGMETRICS)*, Phoenix, AZ, USA, June 2019.
[Preliminary arXiv Version]
[Abstract]
[Slides (pptx) (pdf)]
[MemBen Benchmark Suite]
[Source Code for GPGPUSim-Ramulator]

Demystifying Complex Workload–DRAM Interactions: An Experimental Study

Saugata Ghose[†]

Tianshi Li[†]

Nastaran Hajinazar^{‡†}

Damla Senol Cali[†]

Onur Mutlu^{§†}

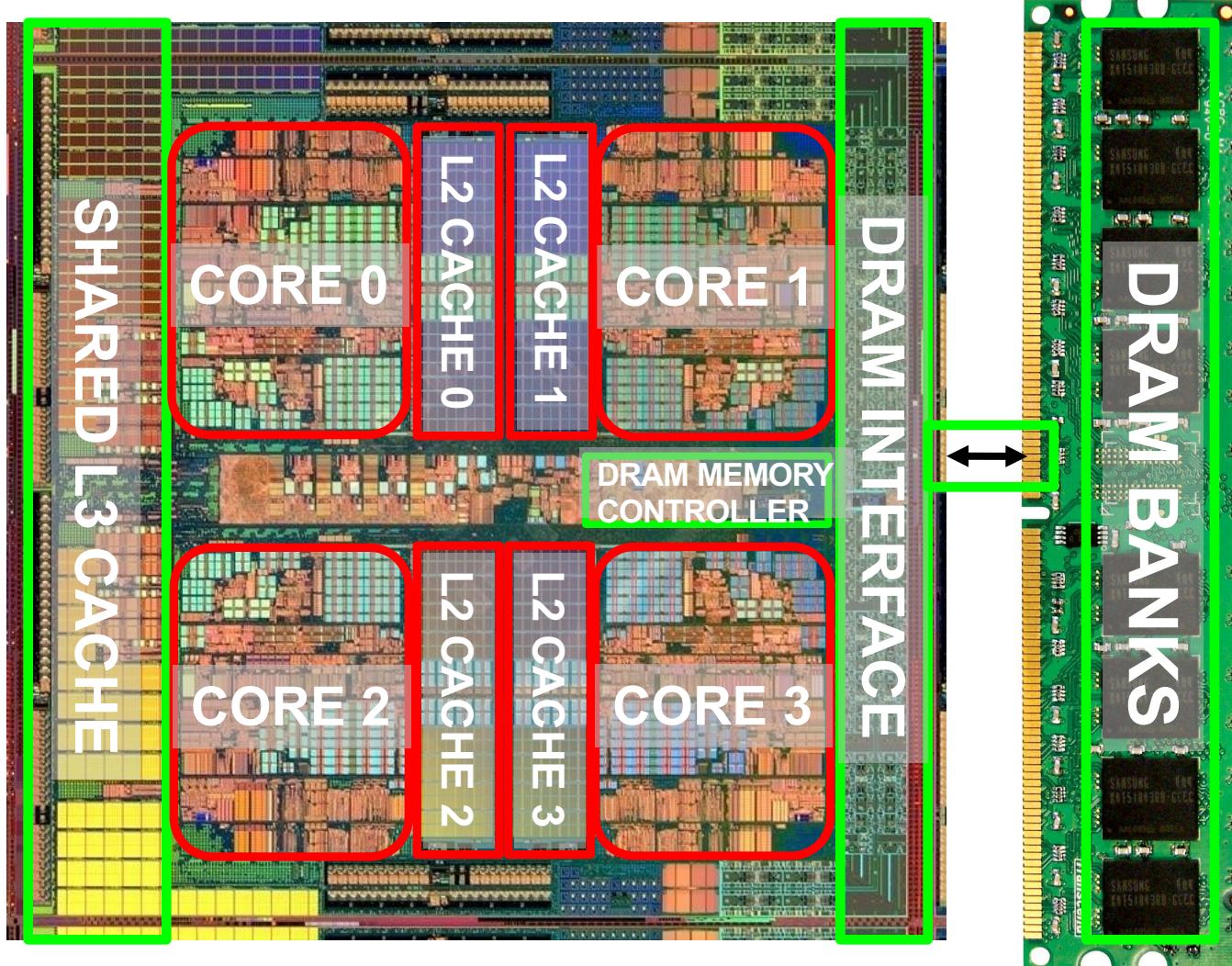
[†]Carnegie Mellon University

[‡]Simon Fraser University

[§]ETH Zürich

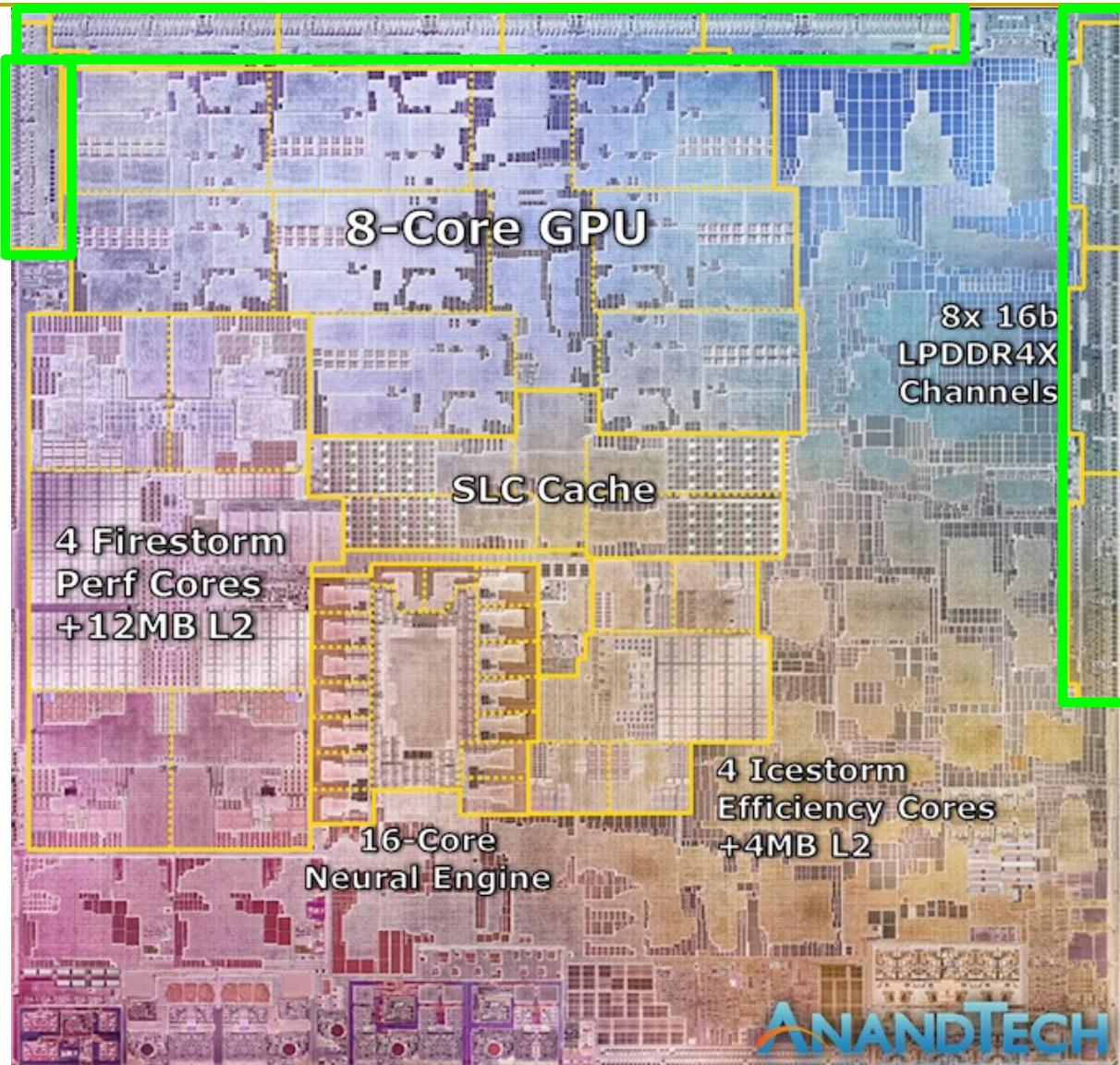
DRAM Control Logic Is Large

Multi-Core
Chip



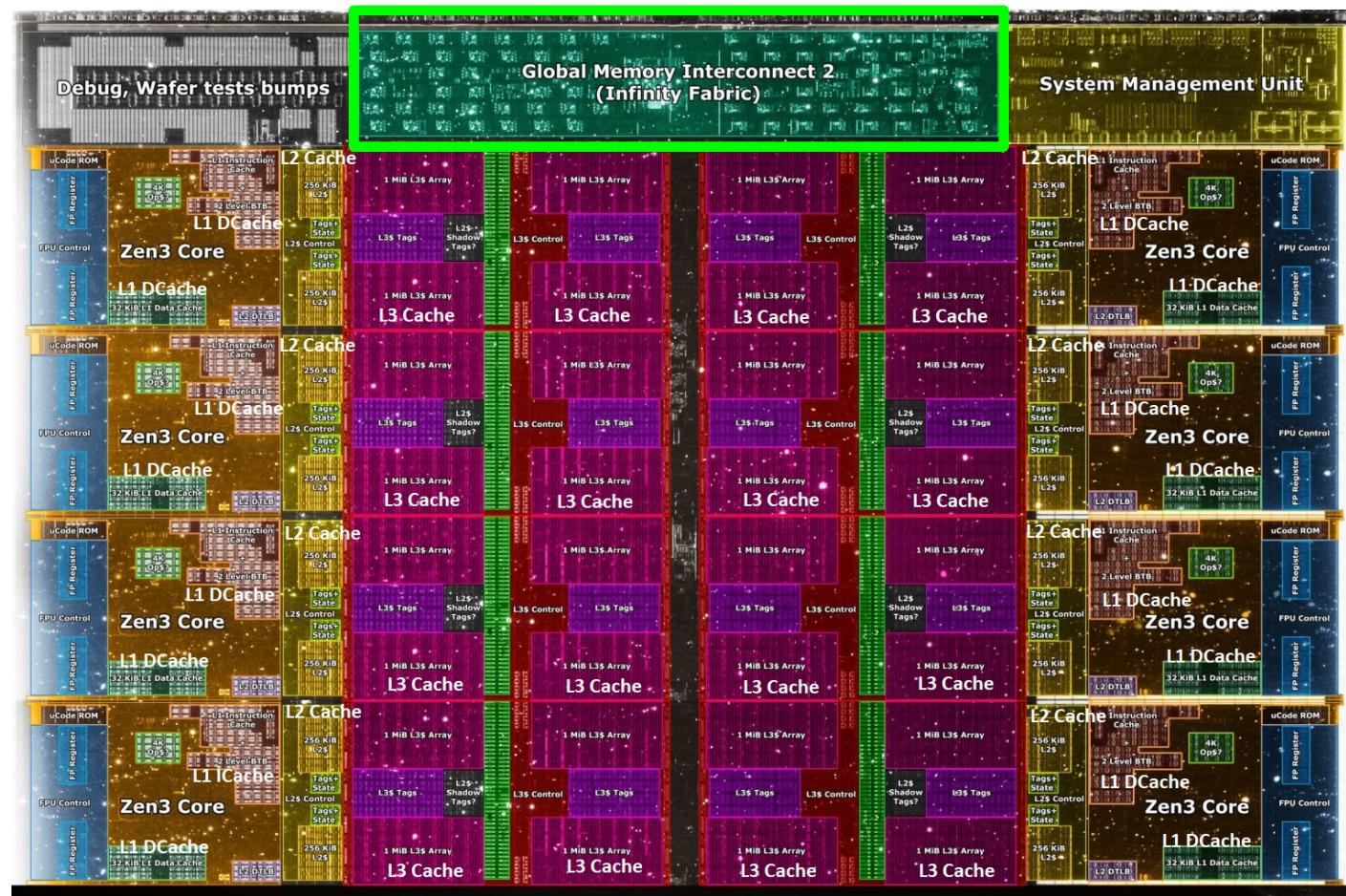
*Die photo credit: AMD Barcelona

DRAM Control Logic Is Large



Apple M1,
2021

DRAM Control Logic Is Large



Core Count:
8 cores/16 threads

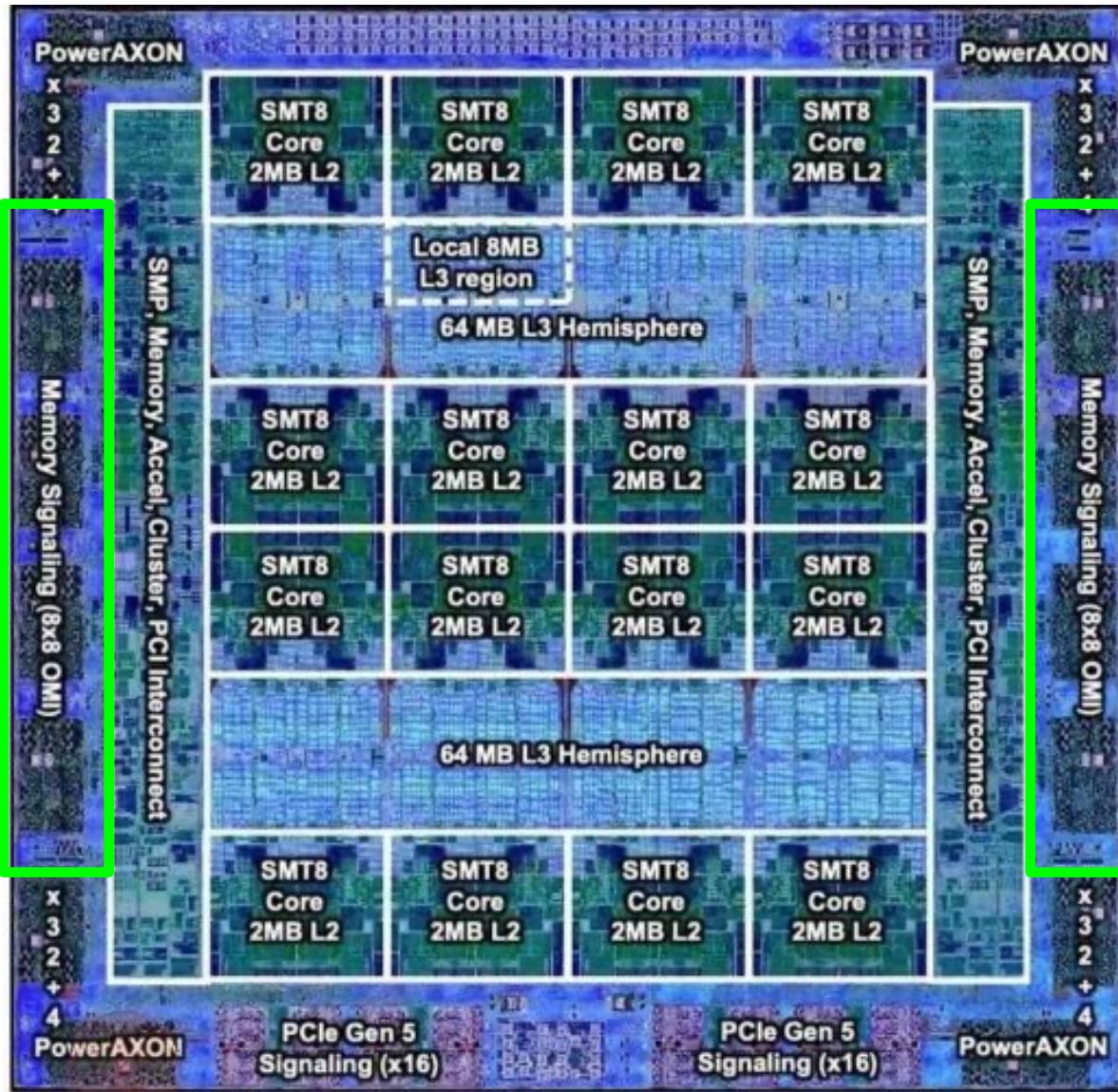
L1 Caches:
32 KB per core

L2 Caches:
512 KB per core

L3 Cache:
32 MB shared

AMD Ryzen 5000, 2020

DRAM Control Logic Is Large



IBM POWER10,
2020

Cores:
15-16 cores,
8 threads/core

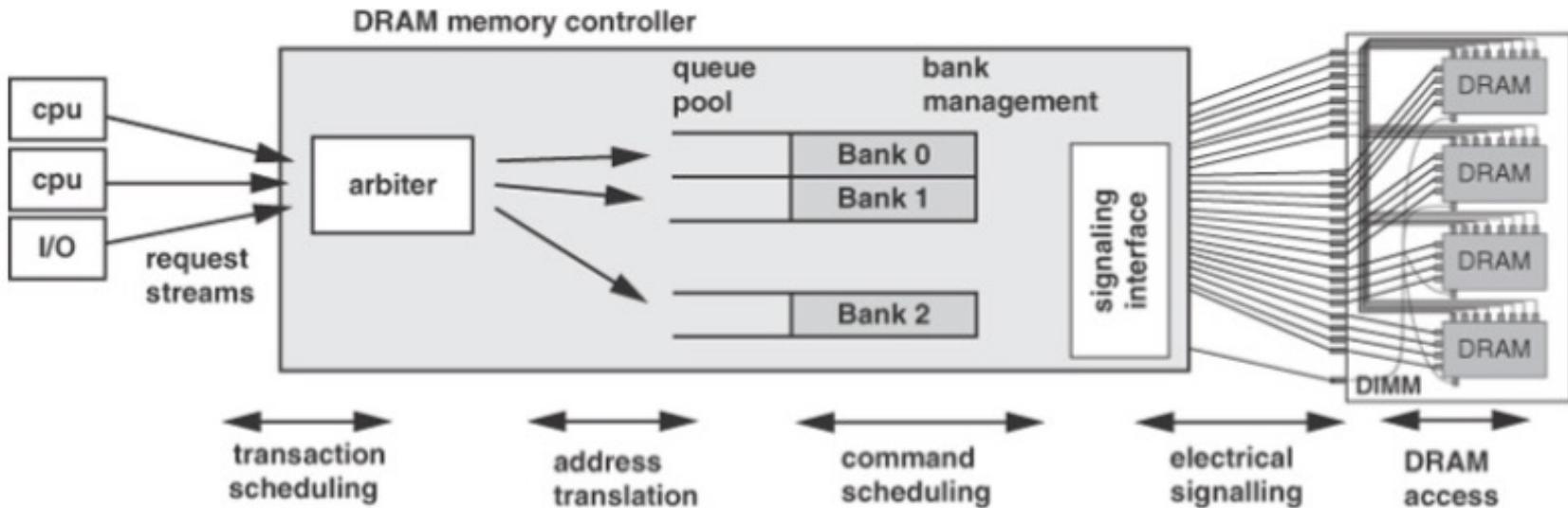
L2 Caches:
2 MB per core

L3 Cache:
120 MB shared

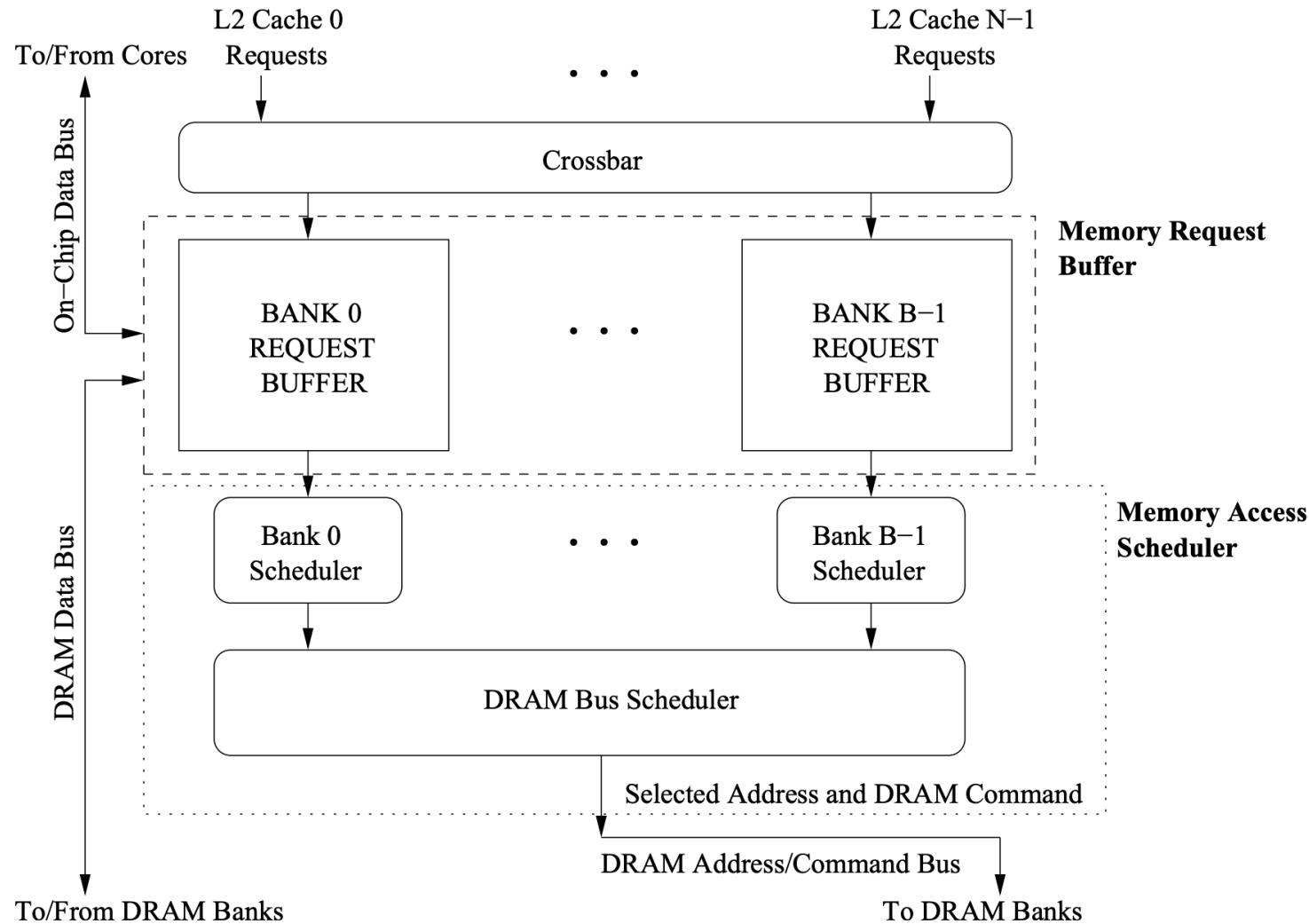
DRAM Controller: Functions

- Ensure correct operation of DRAM (refresh and timing)
- Service DRAM requests while obeying timing constraints of DRAM chips
 - Constraints: resource conflicts (bank, bus, channel), minimum write-to-read delays
 - Translate requests to DRAM command sequences
- Buffer and schedule requests for high performance + QoS
 - Reordering, row-buffer, bank, rank, bus management
- Manage power consumption and thermals in DRAM
 - Turn on/off DRAM chips, manage power modes

A Modern DRAM Controller (I)



A Modern DRAM Controller



DRAM Scheduling Policies (I)

- FCFS (first come first served)
 - Oldest request first
 - FR-FCFS (first ready, first come first served)
 1. Row-hit first
 2. Oldest first
- Goal: Maximize row buffer hit rate → maximize DRAM throughput
- Actually, scheduling is done at the command level
 - Column commands (read/write) prioritized over row commands (activate/precharge)
 - Within each group (e.g., bank), older commands prioritized over younger ones

DRAM Bank Operation

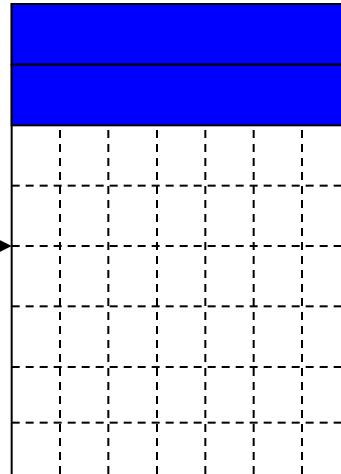
Access Address:

- (Row 0, Column 0)
- (Row 0, Column 1)
- (Row 0, Column 85)
- (Row 1, Column 0)

Row address 0



Columns

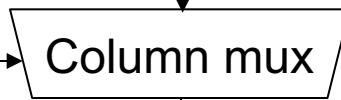


Rows

Row 1

Row Buffer ~~CONF~~ CONFLICT !

Column address 05



Data

DRAM Scheduling Policies (II)

- A scheduling policy is a request prioritization order
- Prioritization can be based on
 - Request age
 - Row buffer hit/miss status
 - Request type (prefetch, read, write)
 - Requestor type (load miss or store miss)
 - Request criticality
 - Oldest miss in the core?
 - How many instructions in core are dependent on it?
 - Will it stall the processor?
 - Interference caused to other cores
 - ...

Row Buffer Management Policies

- Open row
 - Keep the row open after an access
 - + Next access might need the same row → row hit
 - Next access might need a different row → row conflict, wasted energy
- Closed row
 - Close the row after an access (if no other requests already in the request buffer need the same row)
 - + Next access might need a different row → avoid a row conflict
 - Next access might need the same row → extra activate latency
- Adaptive policies
 - Predict whether or not the next access to the bank will be to the same row and act accordingly

Open vs. Closed Row Policies

Policy	First access	Next access	Commands needed for next access
Open row	Row 0	Row 0 (row hit)	Read
Open row	Row 0	Row 1 (row conflict)	Precharge + Activate Row 1 + Read
Closed row	Row 0	Row 0 – access in request buffer (row hit)	Read
Closed row	Row 0	Row 0 – access not in request buffer (row closed)	Activate Row 0 + Read + Precharge
Closed row	Row 0	Row 1 (row closed)	Activate Row 1 + Read + Precharge

DRAM Power Management

- DRAM chips have power modes
- Idea: When not accessing a chip power it down

- Power states
 - Active (highest power)
 - All banks idle
 - Power-down
 - Self-refresh (lowest power)

- Tradeoff: State transitions incur latency during which the chip cannot be accessed

Difficulty of DRAM Control

Why Are DRAM Controllers Difficult to Design?

- Need to obey **DRAM timing constraints** for correctness
 - There are many (50+) timing constraints in DRAM
 - tWTR: Minimum number of cycles to wait before issuing a read command after a write command is issued
 - tRC: Minimum number of cycles between the issuing of two consecutive activate commands to the same bank
 - ...
- Need to **keep track of many resources** to prevent conflicts
 - Channels, banks, ranks, data bus, address bus, row buffers
- Need to handle **DRAM refresh** and **RowHammer**
- Need to **manage power** consumption
- Need to **optimize performance & QoS** (in the presence of constraints)
 - Reordering is not simple
 - Fairness and QoS needs complicates the scheduling problem

Many DRAM Timing Constraints

Latency	Symbol	DRAM cycles	Latency	Symbol	DRAM cycles
Precharge	t_{RP}	11	Activate to read/write	t_{RCD}	11
Read column address strobe	CL	11	Write column address strobe	CWL	8
Additive	AL	0	Activate to activate	t_{RC}	39
Activate to precharge	t_{RAS}	28	Read to precharge	t_{RTP}	6
Burst length	t_{BL}	4	Column address strobe to column address strobe	t_{CCD}	4
Activate to activate (different bank)	t_{RRD}	6	Four activate windows	t_{FAW}	24
Write to read	t_{WTR}	6	Write recovery	t_{WR}	12

Table 4. DDR3 1600 DRAM timing specifications

- From Lee et al., “DRAM-Aware Last-Level Cache Writeback: Reducing Write-Caused Interference in Memory Systems,” HPS Technical Report, April 2010.

More on DRAM Operation

- Kim et al., “A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM,” ISCA 2012.
- Lee et al., “Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture,” HPCA 2013.

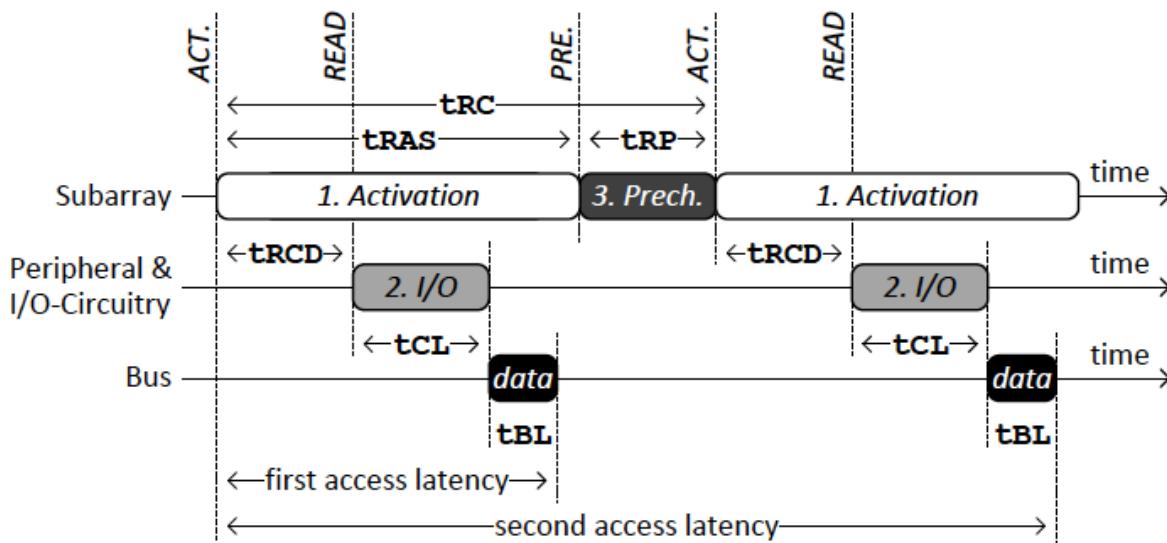
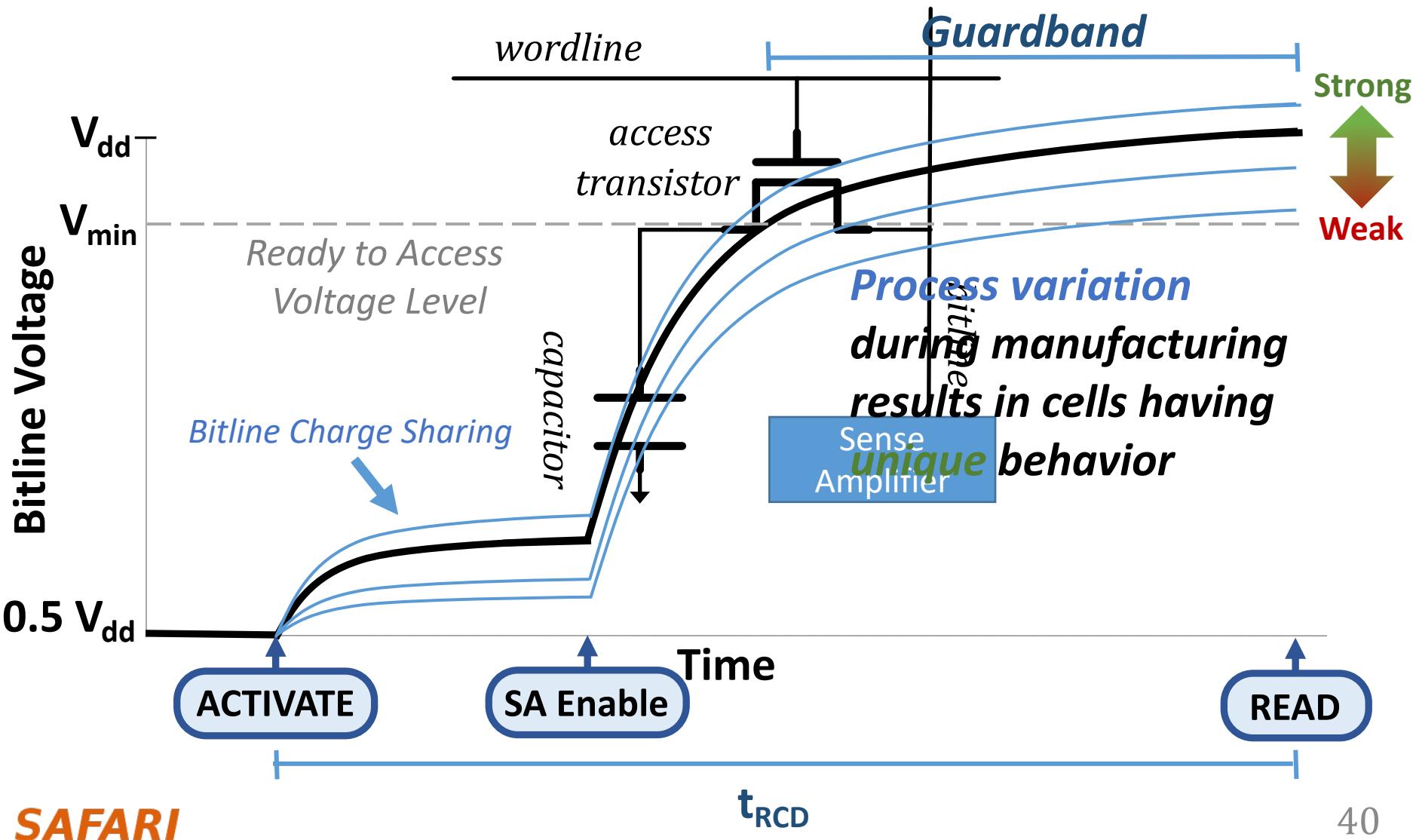


Figure 5. Three Phases of DRAM Access

Table 2. Timing Constraints (DDR3-1066) [43]

Phase	Commands	Name	Value
1	ACT → READ	tRCD	15ns
	ACT → WRITE	tRAS	37.5ns
2	ACT → PRE	tRP	15ns
	READ → data	tCWL	11.25ns
	data burst	tBL	7.5ns
3	PRE → ACT	tRP	15ns
1 & 3	ACT → ACT	tRC (tRAS+tRP)	52.5ns

Why Timing Constraints?



Why So Many Timing Constraints? (I)

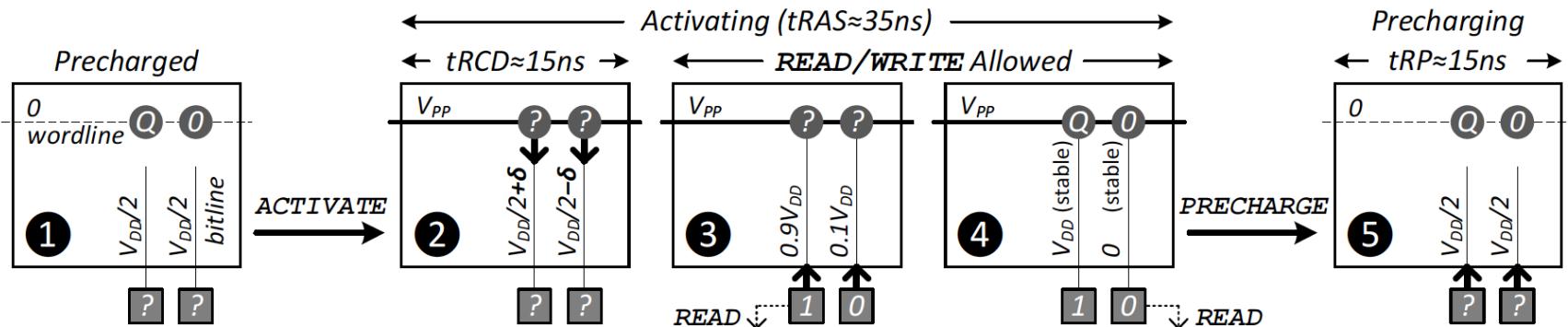


Figure 4. DRAM bank operation: Steps involved in serving a memory request [17] ($V_{PP} > V_{DD}$)

Category	RowCmd \leftrightarrow RowCmd			RowCmd \leftrightarrow ColCmd			ColCmd \leftrightarrow ColCmd			ColCmd \rightarrow DATA	
Name	tRC	$tRAS$	tRP	$tRCD$	$tRTP$	tWR^*	$tCCD$	$tRTW^\dagger$	$tWTR^*$	CL	CWL
Commands	A \rightarrow A	A \rightarrow P	P \rightarrow A	A \rightarrow R/W	R \rightarrow P	W* \rightarrow P	R(W) \rightarrow R(W)	R \rightarrow W	W* \rightarrow R	R \rightarrow DATA	W \rightarrow DATA
Scope	Bank	Bank	Bank	Bank	Bank	Bank	Channel	Rank	Rank	Bank	Bank
Value (ns)	~50	~ 35	13-15	13-15	~ 7.5	15	5-7.5	11-15	~ 7.5	13-15	10-15

A: ACTIVATE – P: PRECHARGE – R: READ – W: WRITE

* Goes into effect after the last write *data*, not from the WRITE command

† Not explicitly specified by the JEDEC DDR3 standard [18]. Defined as a function of other timing constraints.

Table 1. Summary of DDR3-SDRAM timing constraints (derived from Micron's 2Gb DDR3-SDRAM datasheet [33])

Kim et al., "A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM," ISCA 2012.

Why So Many Timing Constraints? (II)

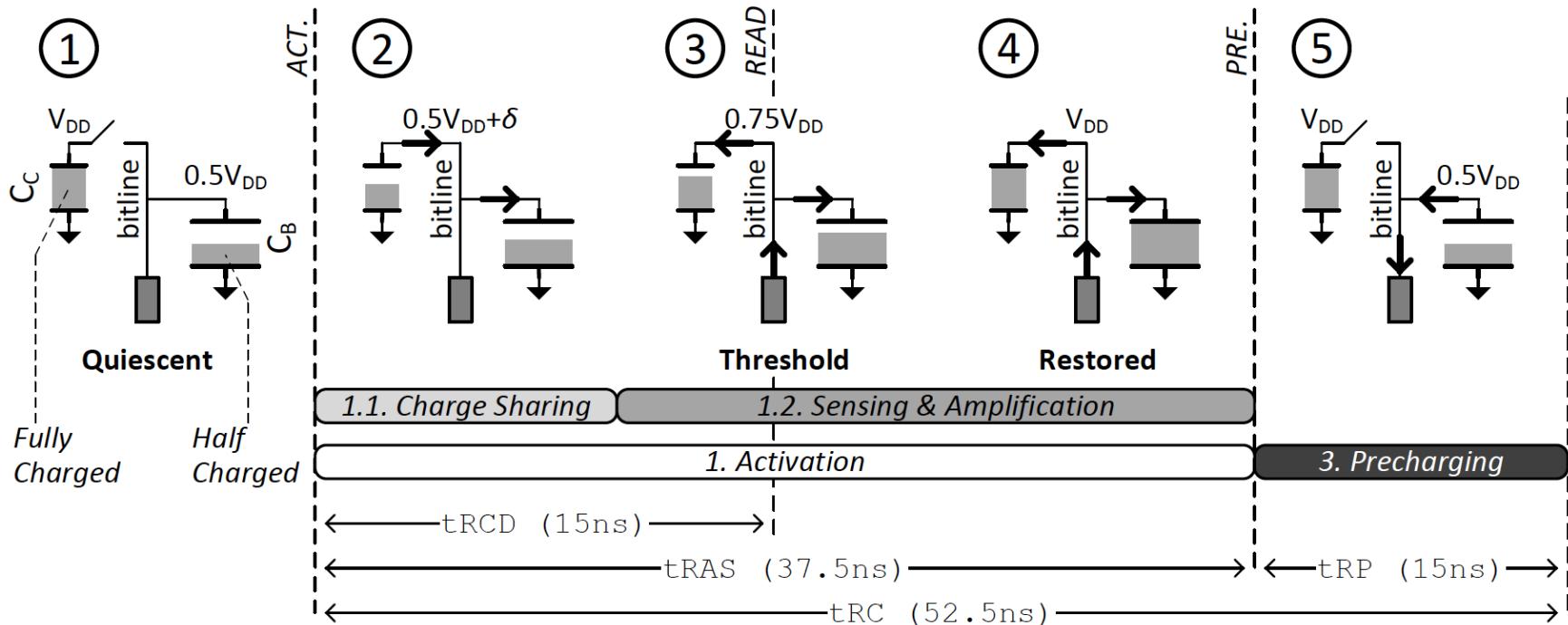


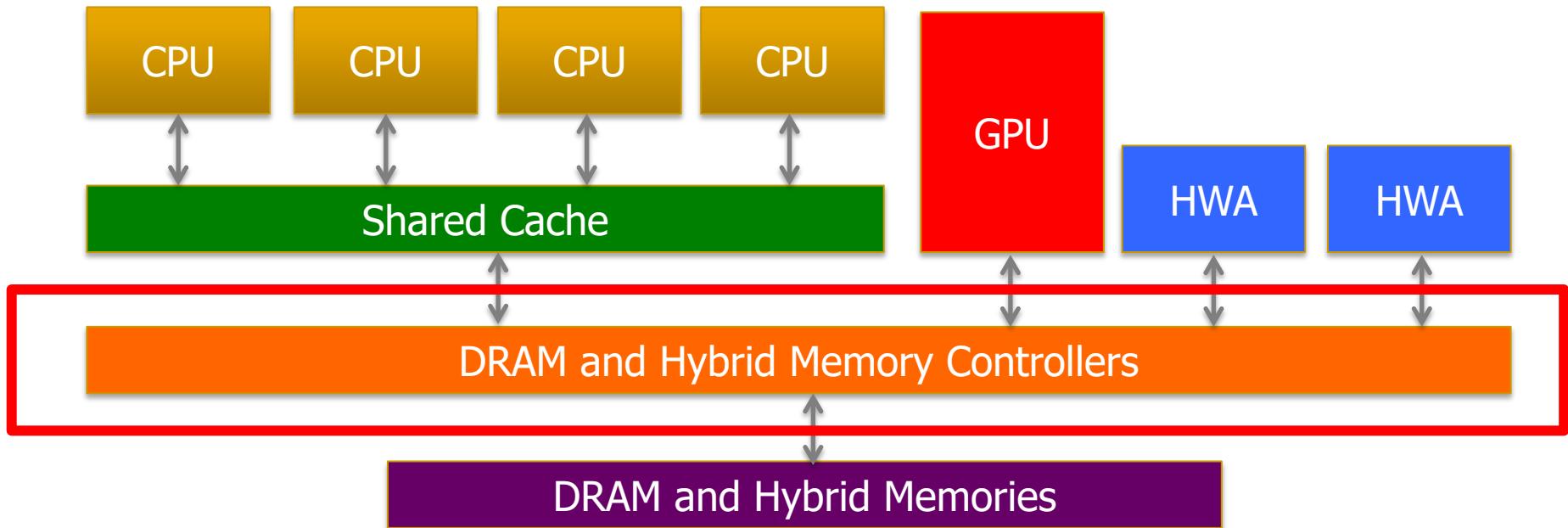
Figure 6. Charge Flow Between the Cell Capacitor (C_C), Bitline Parasitic Capacitor (C_B), and the Sense-Amplifier ($C_B \approx 3.5C_C$ [39])

Lee et al., "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.

Table 2. Timing Constraints (DDR3-1066) [43]

Phase	Commands	Name	Value
1	ACT → READ	t_{RCD}	15ns
	ACT → WRITE	t_{RAS}	37.5ns
2	ACT → PRE	t_{RC}	($t_{RAS} + t_{RP}$) 52.5ns
	READ → <i>data</i>	t_{CL}	15ns
	WRITE → <i>data</i>	t_{CWL}	11.25ns
3	<i>data burst</i>	t_{TBL}	7.5ns
	PRE → ACT	t_{RP}	15ns
1 & 3	ACT → ACT	t_{RC}	52.5ns

DRAM Controller Design Is Becoming More Difficult



- Heterogeneous agents: CPUs, GPUs, and HWAs
- Main memory interference between CPUs, GPUs, HWAs
- Many timing constraints for various memory types
- Many goals at the same time: performance, fairness, QoS, energy efficiency, ...

Self-Managing DRAM

- Hasan Hassan, Ataberk Olgun, A Giray Yaglikci,
Haocong Luo, and Onur Mutlu,

"A Case for Self-Managing DRAM Chips: Improving Performance, Efficiency, Reliability, and Security via Autonomous in-DRAM Maintenance Operations"

Preprint on *arxiv*, March 2023.

[arXiv version]

[SMD Source Code]

**A Case for Self-Managing DRAM Chips:
Improving Performance, Efficiency, Reliability, and Security
via Autonomous in-DRAM Maintenance Operations**

Hasan Hassan

Ataberk Olgun

A. Giray Yağlıkçı

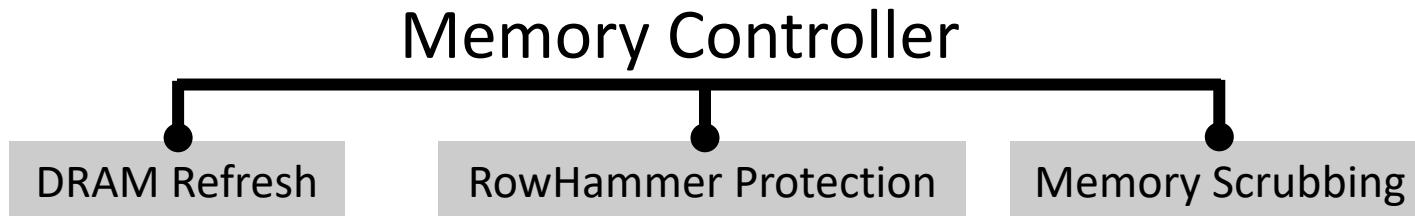
Haocong Luo

Onur Mutlu

ETH Zürich

Self-Managing DRAM: Problem

The **Memory Controller** manages
DRAM **maintenance operations**



Changes to **maintenance operations** are **often reflected** to the memory controller design, DRAM interface, and other system components



Implementing new maintenance operations
(or modifying the existing ones) is **difficult-to-realize**

A Prime Example: New Features of DDR5

DRAM Refresh

Same Bank Refresh - simultaneously refreshes one bank in each bank group

The new **REFsb** command requires changes in DRAM interface and memory controller

RowHammer Protection

Refresh Management (RFM) – memory controller issues the new RFM command to allow DRAM chips more time to perform victim row refresh

The new **RFM** command requires changes in DRAM interface and memory controller

Memory Scrubbing

In-DRAM Scrubbing – DDR5 uses the on-die ECC to perform periodic scrubbing

The new **scrub** command requires changes in DRAM interface and memory controller

DDR5 changes are difficult-to-implement as they were *only* possible after multiple years required to develop a new DRAM standard

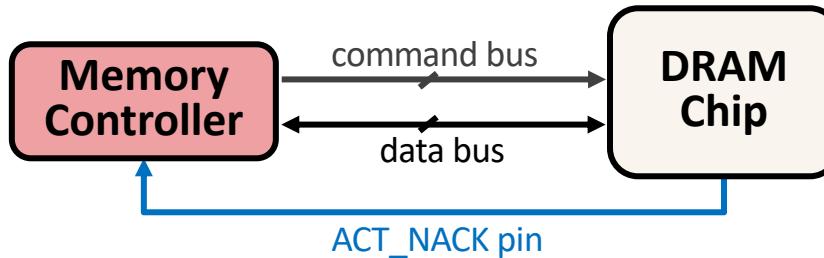
SMD: Overview

Self-Managing DRAM (SMD)

enables autonomous in-DRAM maintenance operations

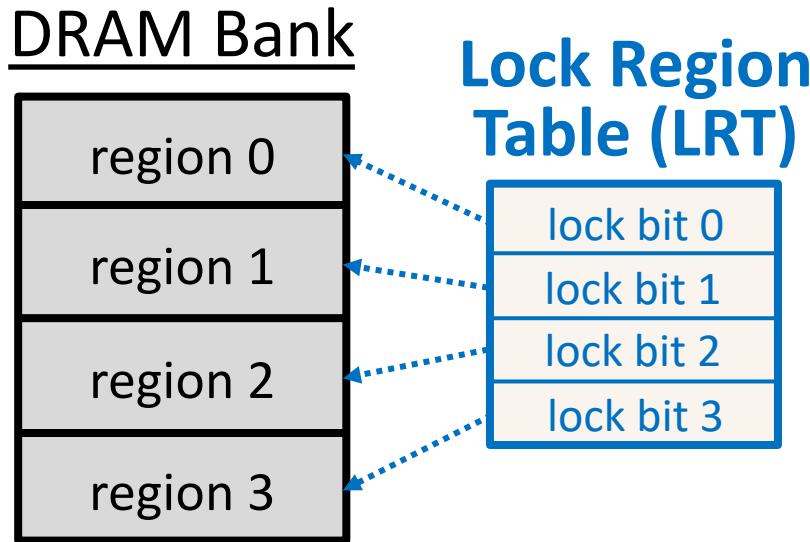
Key Idea:

Prevent the memory controller from accessing DRAM regions that are *under maintenance* by **rejecting** row activation (ACT) commands



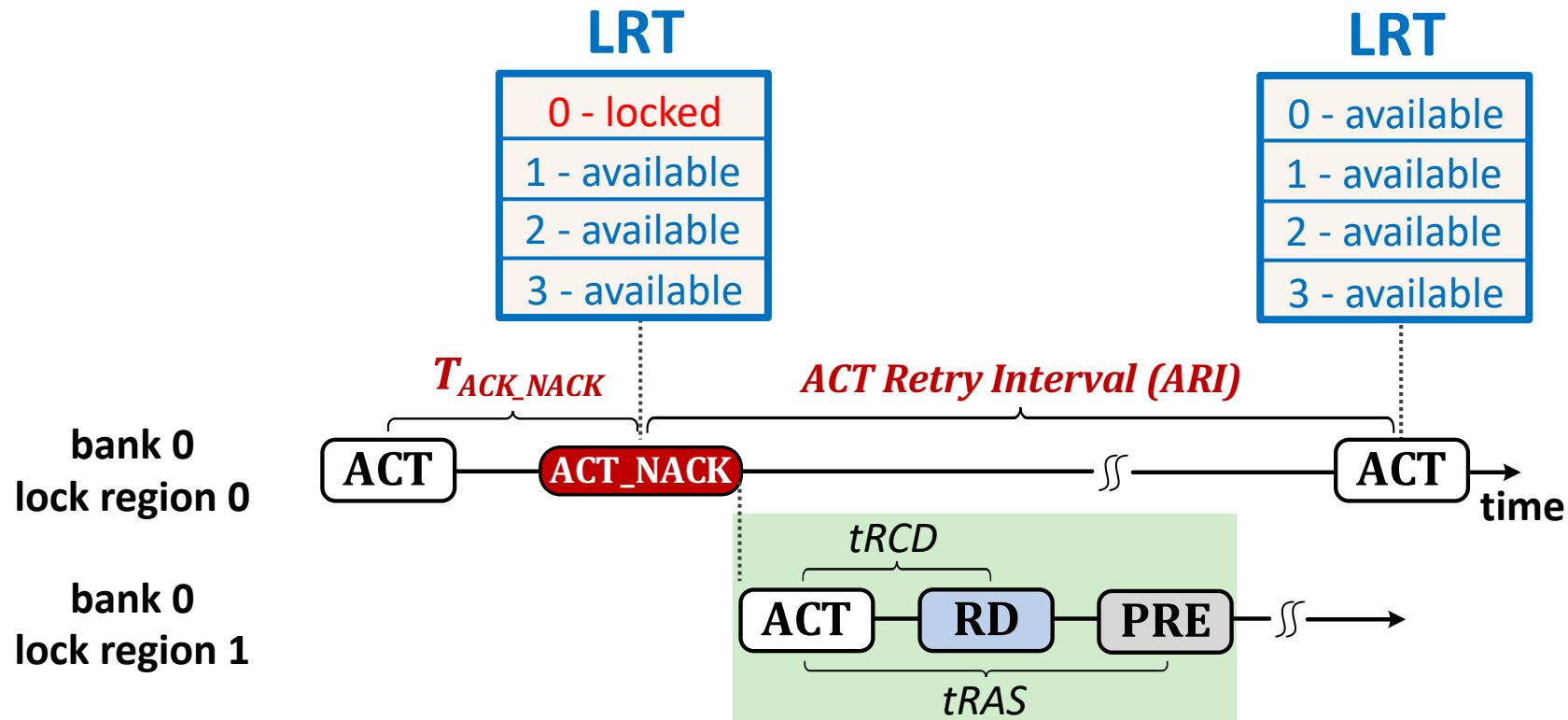
Leveraging the ability to *reject an ACT*, a **maintenance operation** can be implemented **completely** within a DRAM chip

SMD: DRAM Bank Architecture



- A DRAM bank is divided into configurable-size **Lock Regions**
- SMD marks regions *locked* for maintenance in the **Lock Region Table (LRT)**
- SMD **rejects** any ACT command that targets a *locked region* by sending the memory controller an **ACT_NACK** signal

SMD: High-Level Operation



SMD-Based Maintenance Mechanisms

DRAM Refresh

Fixed Rate (SMD-FR)

*uniformly refreshes all DRAM rows with a **fixed** refresh period*

Variable Rate (SMD-VR)

*skips refreshing rows that can **retain their data for longer** than the default refresh period*

RowHammer Protection

Probabilistic (SMD-PRP)

*Performs **neighbor row refresh** with a **small probability** on every row activation*

Deterministic (SMD-DRP)

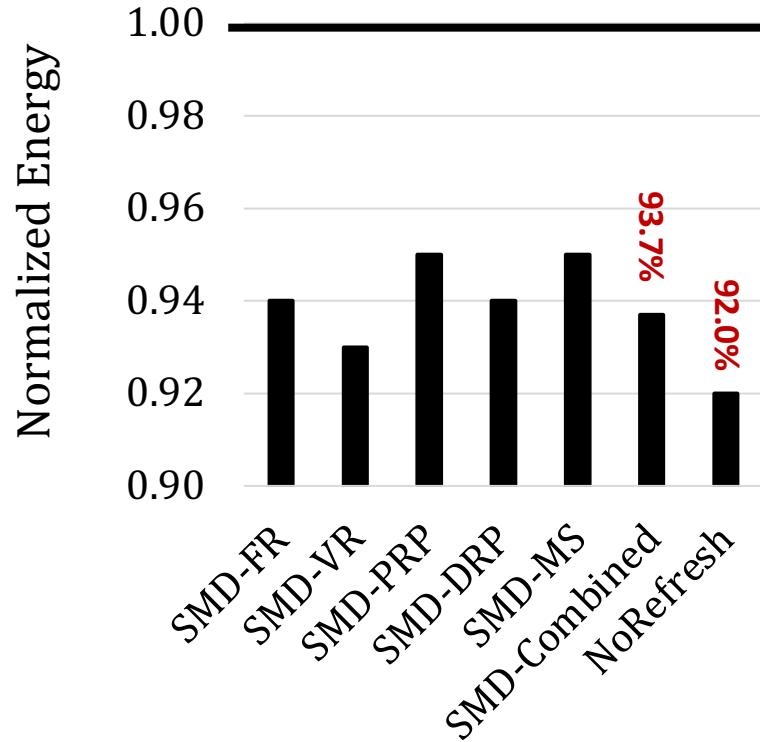
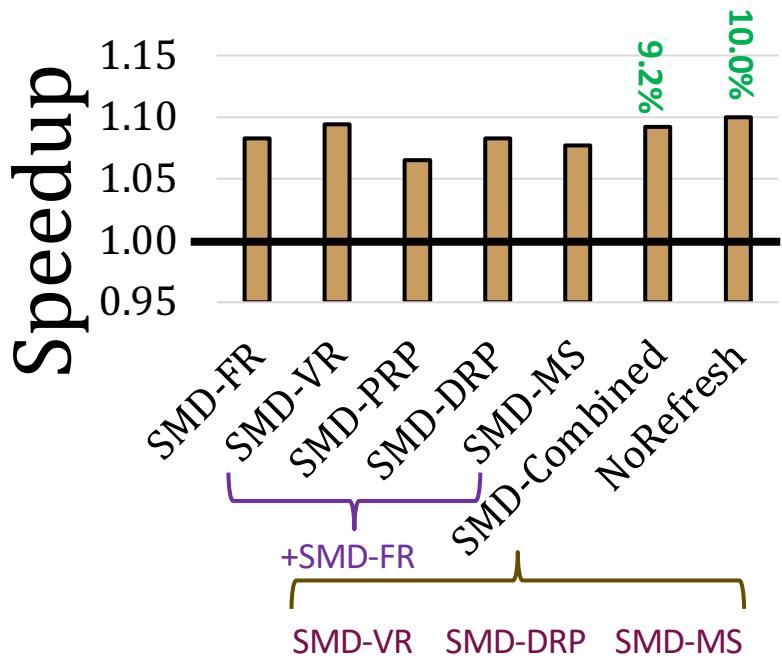
*keeps track of most **frequently activated** rows and performs **neighbor row refresh** when activation count threshold is exceeded*

Memory Scrubbing

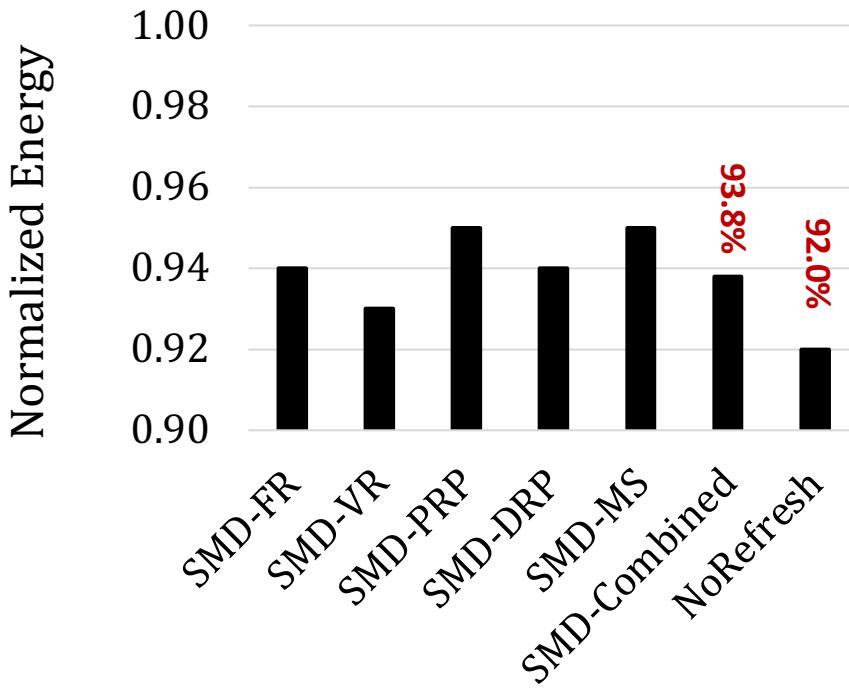
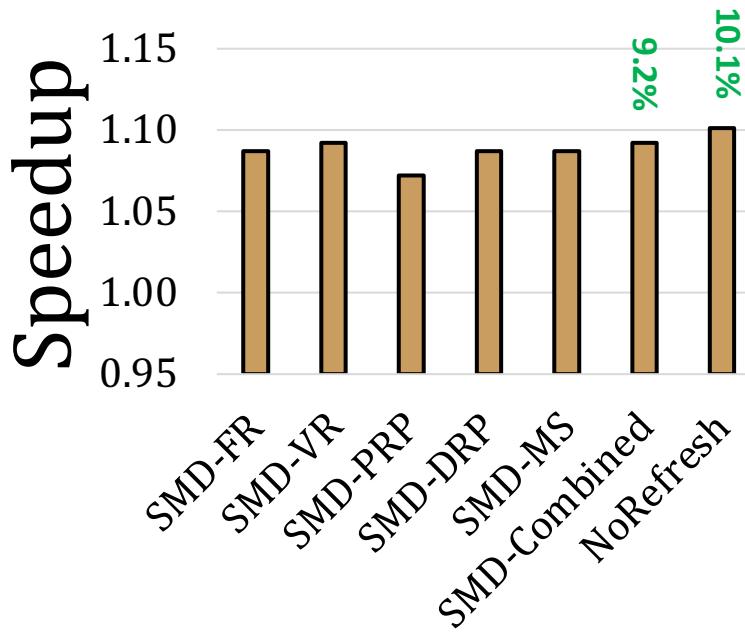
Periodic Scrubbing (SMD-MS)

*periodically **scans** the **entire DRAM** for errors and corrects them*

Single-Core Results

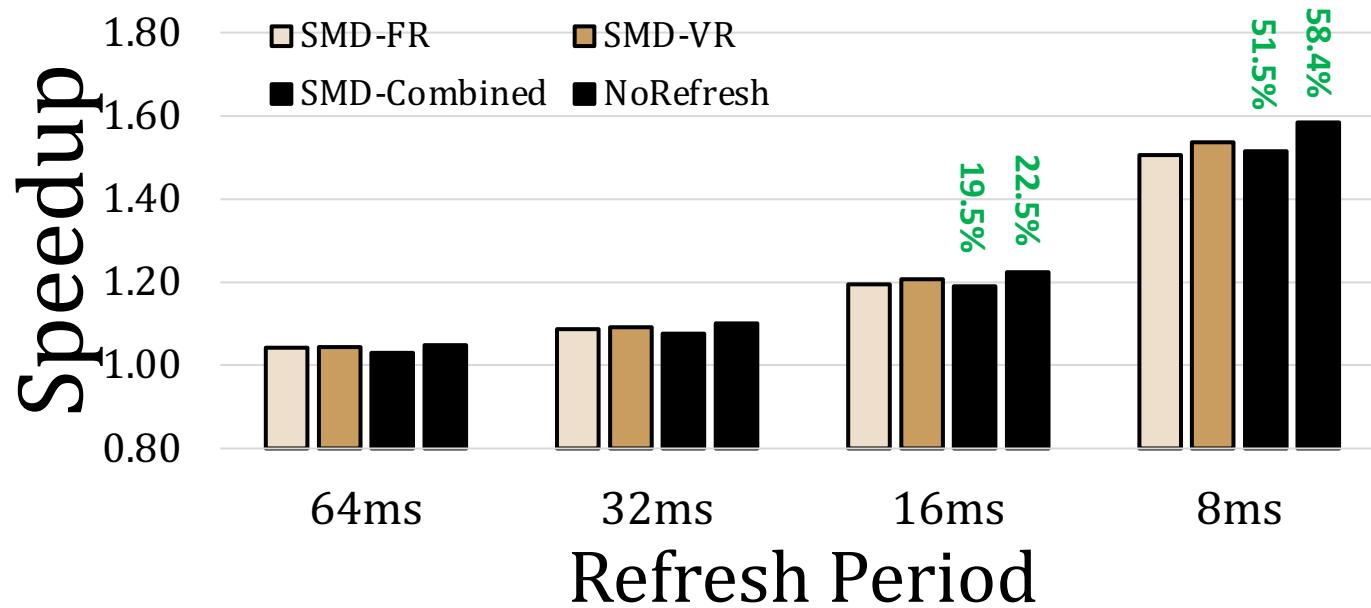


Four-Core Results



SMD-based maintenance mechanisms have **significant**
performance and energy efficiency **benefits**

Sensitivity to Refresh Period



Hardware Overhead

Interface Modifications

- A single **ACT_NACK** pin per DRAM chip or **repurpose** the existing alert_n signal

DRAM Chip Modifications

- **Lock Region Table** incurs only:
 - 32um² **area** overhead (0.001% of a 45.4mm² DRAM chip)
 - 0.053ns **access latency** overhead

Memory Controller Modifications

- **Changes** in *request scheduling* to handle ACT_NACK signals

- **No further changes needed** for new maintenance operations
- The memory controller **no longer manages** DRAM maintenance

Other Results in the Paper

- Lock region size sensitivity
- Comparison to memory controller-based RH protection
- Comparison to memory controller-based scrubbing
- SMD-DRP maximum activation threshold sensitivity
- Victim row window sensitivity

Summary

The three major DRAM maintenance operations:

- ❖ Refresh
- ❖ RowHammer Protection
- ❖ Memory Scrubbing

Source code is available:

github.com/CMU-SAFARI/SelfManagingDRAM



Implementing new **maintenance mechanisms** often requires **difficult-to-realize changes**

Our Goal

- ① Ease the process of enabling new DRAM maintenance operations
- ② Enable more efficient in-DRAM maintenance operations

Self-Managing DRAM (SMD)

Enables implementing new **in-DRAM** maintenance mechanisms
with **no further changes** in the *DRAM interface* and *memory controller*

SMD-based *refresh*, *RowHammer protection*, and *scrubbing* achieve
9.2% speedup and **6.2% lower DRAM energy** vs. conventional DRAM

Self-Managing DRAM

- Hasan Hassan, Ataberk Olgun, A Giray Yaglikci,
Haocong Luo, and Onur Mutlu,

"A Case for Self-Managing DRAM Chips: Improving Performance, Efficiency, Reliability, and Security via Autonomous in-DRAM Maintenance Operations"

Preprint on *arxiv*, March 2023.

[arXiv version]

[SMD Source Code]

**A Case for Self-Managing DRAM Chips:
Improving Performance, Efficiency, Reliability, and Security
via Autonomous in-DRAM Maintenance Operations**

Hasan Hassan

Ataberk Olgun

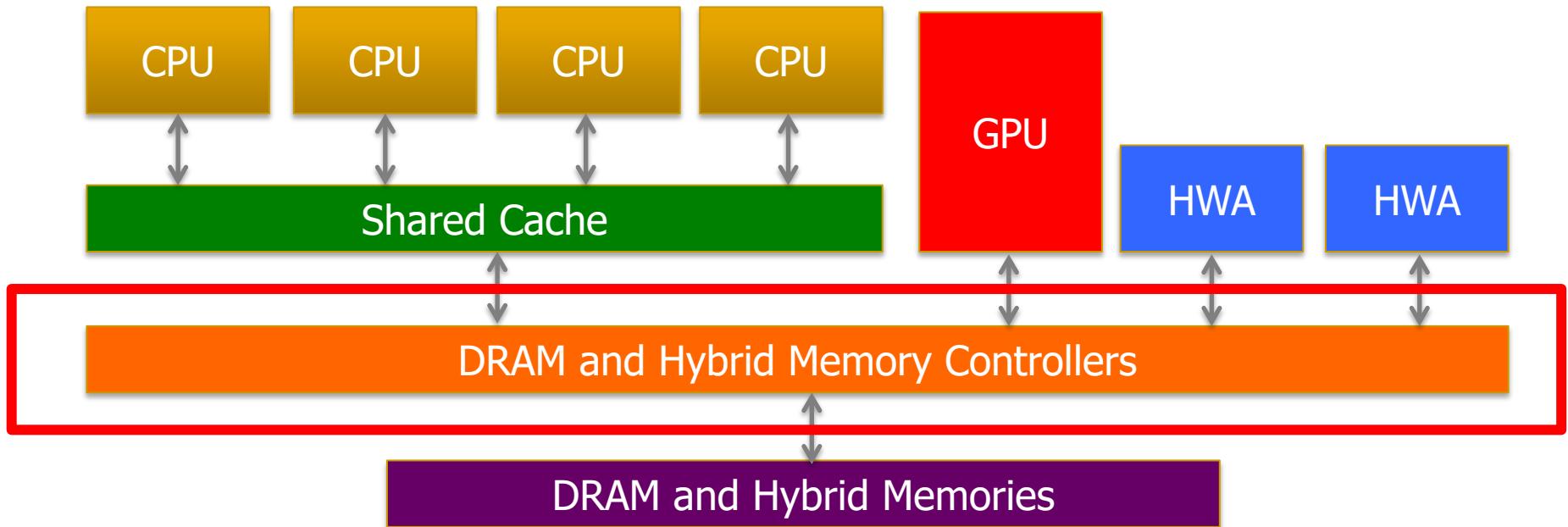
A. Giray Yağlıkçı

Haocong Luo

Onur Mutlu

ETH Zürich

DRAM Controller Design Is Becoming More Difficult

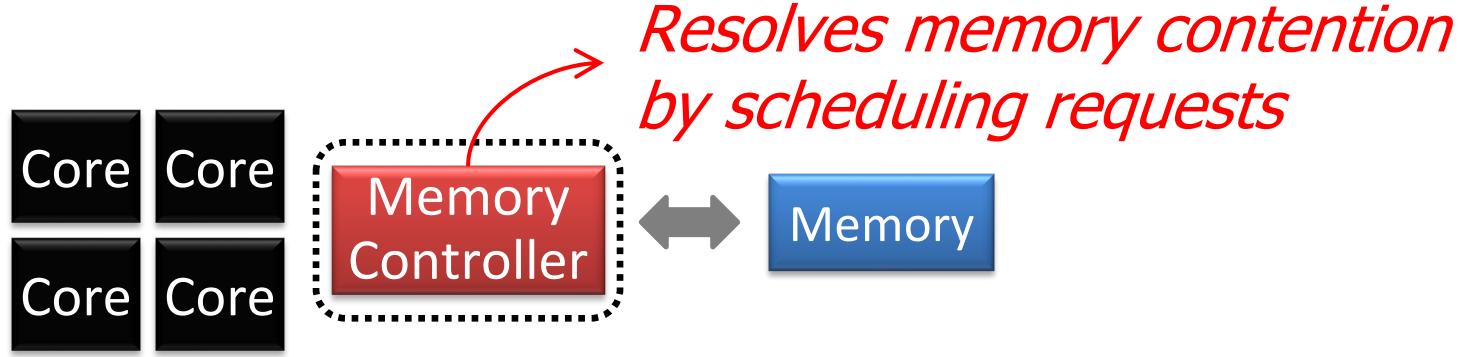


- Heterogeneous agents: CPUs, GPUs, and HWAs
- Main memory interference between CPUs, GPUs, HWAs
- Many timing constraints for various memory types
- Many goals at the same time: performance, fairness, QoS, energy efficiency, ...

Reality and Dream

- Reality: It is difficult to design a policy that maximizes performance, QoS, energy-efficiency, ...
 - Too many things to think about
 - Continuously changing workload and system behavior
- Dream: Wouldn't it be nice if the DRAM controller automatically found a good scheduling policy on its own?

Memory Controller: Performance Function

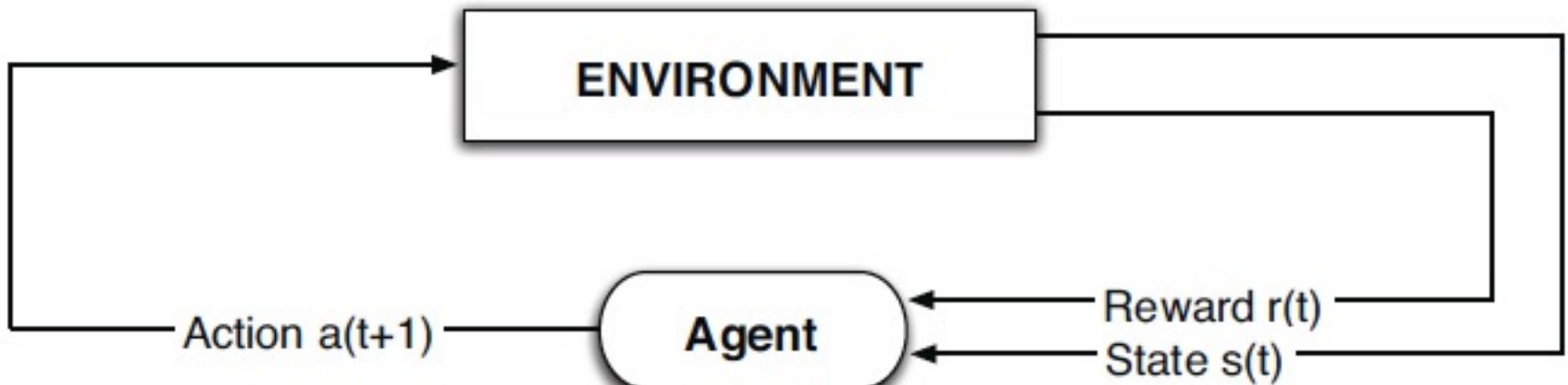


How to schedule requests to maximize system performance?

Self-Optimizing DRAM Controllers

- Problem: DRAM controllers are difficult to design
 - It is difficult for human designers to design a policy that can adapt itself very well to different workloads and different system conditions
- Idea: A memory controller that adapts its scheduling policy to workload behavior and system conditions using machine learning.
- Observation: Reinforcement learning maps nicely to memory control.
- Design: Memory controller is a reinforcement learning agent
 - It dynamically and continuously learns and employs the best scheduling policy to maximize long-term performance.

Self-Optimizing DRAM Controllers

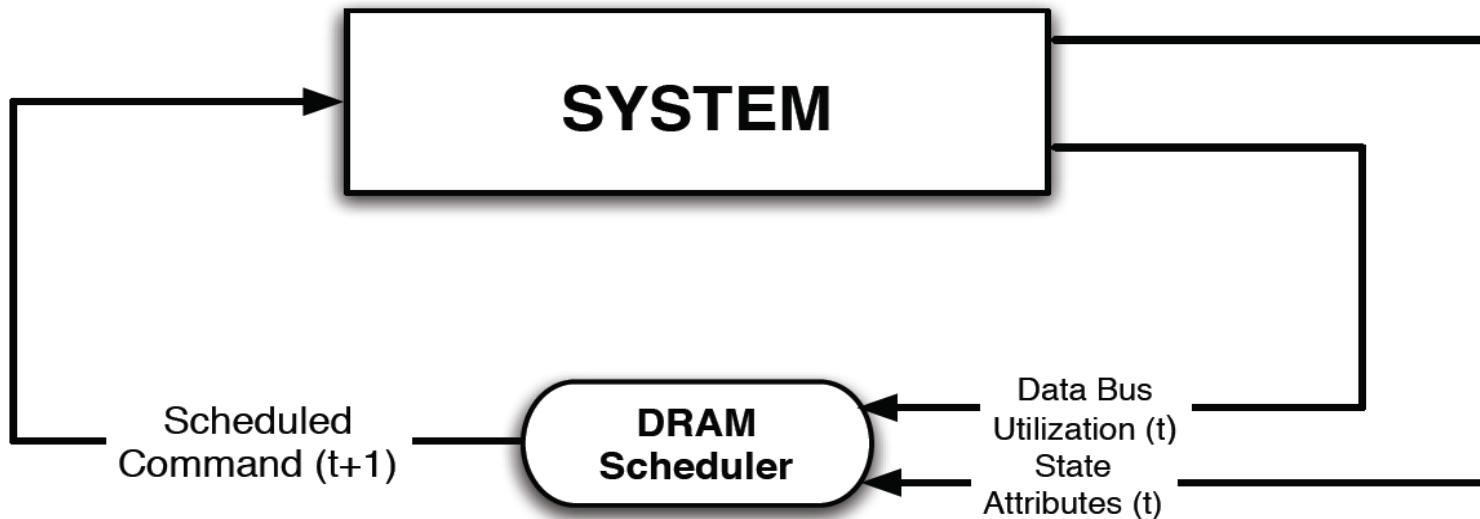


Goal: Learn to choose actions to maximize $r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$ ($0 \leq \gamma < 1$)

Figure 2: (a) Intelligent agent based on reinforcement learning principles;

Self-Optimizing DRAM Controllers

- Dynamically adapt the memory scheduling policy via interaction with the system at runtime
 - Associate system states and actions (commands) with long term reward values: each action at a given state leads to a learned reward
 - Schedule command with highest estimated long-term reward value in each state
 - Continuously update reward values for <state, action> pairs based on feedback from system



Self-Optimizing DRAM Controllers

- 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.

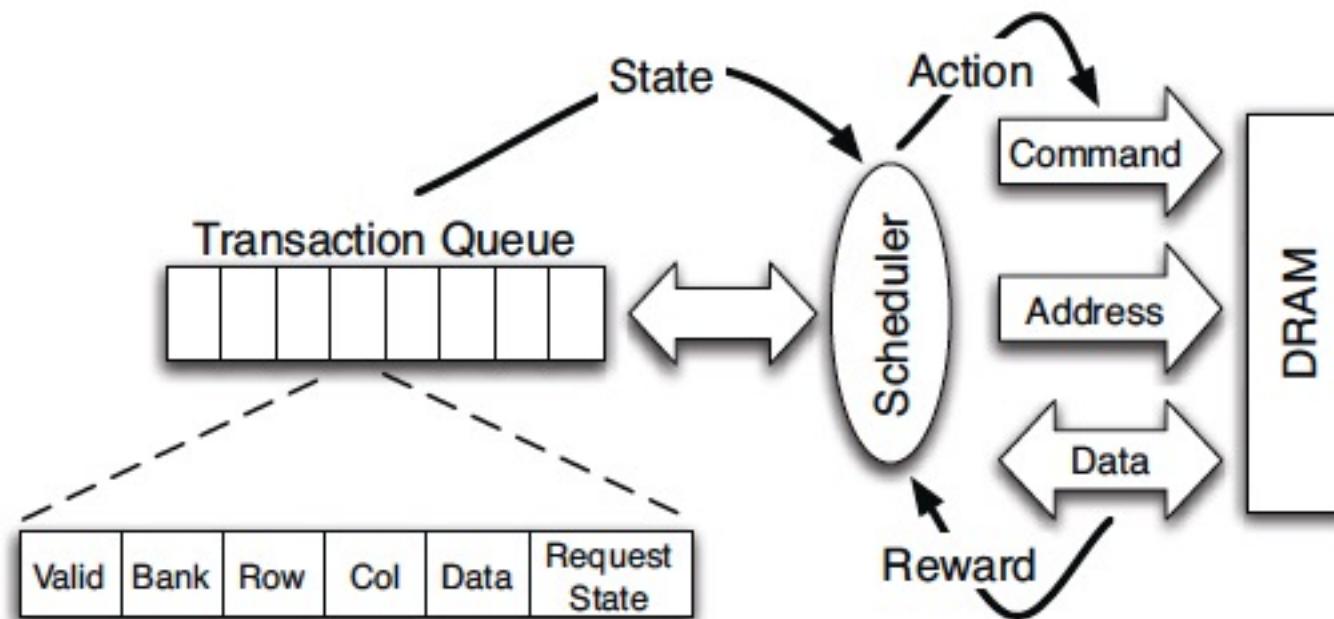


Figure 4: High-level overview of an RL-based scheduler.

States, Actions, Rewards

❖ Reward function

- +1 for scheduling Read and Write commands
- 0 at all other times

Goal is to maximize long-term data bus utilization

❖ State attributes

- Number of reads, writes, and load misses in transaction queue
- Number of pending writes and ROB heads waiting for referenced row
- Request's relative ROB order

❖ Actions

- Activate
- Write
- Read - load miss
- Read - store miss
- Precharge - pending
- Precharge - preemptive
- NOP

Performance Results

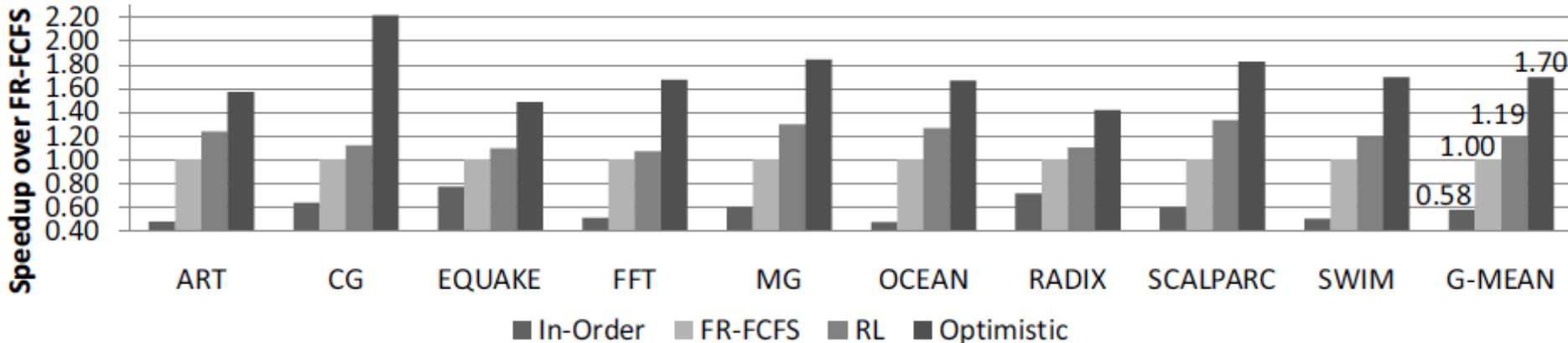


Figure 7: Performance comparison of in-order, FR-FCFS, RL-based, and optimistic memory controllers

Large, robust performance improvements over many human-designed policies

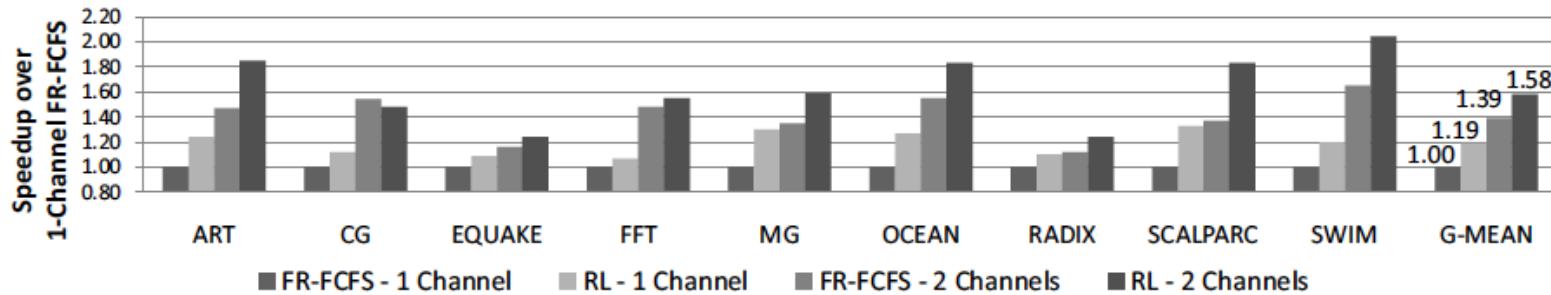


Figure 15: Performance comparison of FR-FCFS and RL-based memory controllers on systems with 6.4GB/s and 12.8GB/s peak DRAM bandwidth

Self Optimizing DRAM Controllers

- + Continuous learning in the presence of changing environment
 - + Reduced designer burden in finding a good scheduling policy.
Designer specifies:
 - 1) What system variables might be useful
 - 2) What target to optimize, but not how to optimize it
 - How to specify different objectives? (e.g., fairness, QoS, ...)
 - Hardware complexity?
 - Design mindset and flow
-

More on Self-Optimizing DRAM Controllers

- 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.

Self-Optimizing Memory Controllers: A Reinforcement Learning Approach

Engin İpek^{1,2} Onur Mutlu² José F. Martínez¹ Rich Caruana¹

¹Cornell University, Ithaca, NY 14850 USA

² Microsoft Research, Redmond, WA 98052 USA

Self-Optimizing (Data-Driven)

Computing Architectures

System Architecture Design Today

- Human-driven
 - Humans design the policies (how to do things)
- Many (too) simple, short-sighted policies all over the system
- No automatic data-driven policy learning
- (Almost) no learning: cannot take lessons from past actions

**Can we design
fundamentally intelligent architectures?**

An Intelligent Architecture

- Data-driven
 - Machine learns the “best” policies (how to do things)
- Sophisticated, workload-driven, changing, far-sighted policies
- Automatic data-driven policy learning
- All controllers are intelligent data-driven agents

**We need to rethink design
(of all controllers)**

Self-Optimizing Memory Prefetchers

Rahul Bera, Konstantinos Kanellopoulos, Anant Nori, Taha Shahroodi, Sreenivas Subramoney, and Onur Mutlu,

"Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning"

Proceedings of the *54th International Symposium on Microarchitecture (MICRO)*, Virtual, October 2021.

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

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

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

[[Talk Video](#) (20 minutes)]

[[Lightning Talk Video](#) (1.5 minutes)]

[[Pythia Source Code \(Officially Artifact Evaluated with All Badges\)](#)]

[[arXiv version](#)]

Officially artifact evaluated as available, reusable and reproducible.



Pythia: A Customizable Hardware Prefetching Framework Using Online Reinforcement Learning

Rahul Bera¹

Konstantinos Kanellopoulos¹

Anant V. Nori²

Taha Shahroodi^{3,1}

Sreenivas Subramoney²

Onur Mutlu¹

¹ETH Zürich

²Processor Architecture Research Labs, Intel Labs

³TU Delft

Learning-Based Off-Chip Load Predictors

- Rahul Bera, Konstantinos Kanellopoulos, Shankar Balachandran, David Novo, Ataberk Olgun, Mohammad Sadrosadati, and Onur Mutlu,
"Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction"

Proceedings of the 55th International Symposium on Microarchitecture (MICRO), Chicago, IL, USA, October 2022.

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

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

[[Talk Video](#) (12 minutes)]

[[Lecture Video](#) (25 minutes)]

[[arXiv version](#)]

[[Source Code \(Officially Artifact Evaluated with All Badges\)](#)]

Officially artifact evaluated as available, reusable and reproducible.

Best paper award at MICRO 2022.



Hermes: Accelerating Long-Latency Load Requests via Perceptron-Based Off-Chip Load Prediction

Rahul Bera¹ Konstantinos Kanellopoulos¹ Shankar Balachandran² David Novo³
Ataberk Olgun¹ Mohammad Sadrosadati¹ Onur Mutlu¹

¹ETH Zürich ²Intel Processor Architecture Research Lab ³LIRMM, Univ. Montpellier, CNRS

Self-Optimizing Hybrid SSD Controllers

Gagandeep Singh, Rakesh Nadig, Jisung Park, Rahul Bera, Nastaran Hajinazar, David Novo, Juan Gomez-Luna, Sander Stuijk, Henk Corporaal, and Onur Mutlu,

"Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning"

Proceedings of the 49th International Symposium on Computer Architecture (ISCA), New York, June 2022.

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

[[arXiv version](#)]

[[Sibyl Source Code](#)]

[[Talk Video](#) (16 minutes)]

Sibyl: Adaptive and Extensible Data Placement in Hybrid Storage Systems Using Online Reinforcement Learning

Gagandeep Singh¹

David Novo³

Rakesh Nadig¹

Juan Gómez-Luna¹

Jisung Park¹

Sander Stuijk²

Rahul Bera¹

Henk Corporaal²

Nastaran Hajinazar¹

Onur Mutlu¹

¹ETH Zürich

²Eindhoven University of Technology

³LIRMM, Univ. Montpellier, CNRS

Architectures for Intelligent Machines

Data-centric

Data-driven

Data-aware

Key Problems with Today's Architectures

- Architectures are **terrible at dealing with data**
 - Designed to mainly store and move data vs. to compute
 - They are **processor-centric** as opposed to **data-centric**
- Architectures are **terrible at taking advantage of vast amounts of data** (and metadata) available to them
 - Designed to make simple decisions, ignoring lots of data
 - They make **human-driven decisions** vs. **data-driven** decisions
- Architectures are **terrible at knowing and exploiting different properties of application data**
 - Designed to treat all data as the same
 - They make **component-aware decisions** vs. **data-aware**

Fundamentally Better Architectures

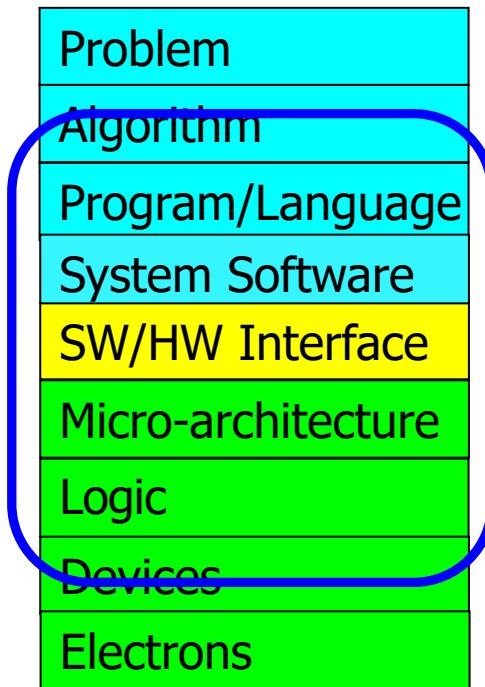
Data-centric

Data-driven

Data-aware



We Need to Think Across the Entire Stack



We can get there step by step

A Blueprint for Fundamentally Better Architectures

- Onur Mutlu,

"Intelligent Architectures for Intelligent Computing Systems"

Invited Paper in Proceedings of the Design, Automation, and Test in Europe Conference (DATE), Virtual, February 2021.

[Slides (pptx) (pdf)]

[IEDM Tutorial Slides (pptx) (pdf)]

[Short DATE Talk Video (11 minutes)]

[Longer IEDM Tutorial Video (1 hr 51 minutes)]

Intelligent Architectures for Intelligent Computing Systems

Onur Mutlu
ETH Zurich
omutlu@gmail.com

A Tutorial on Fundamentally Better Architectures

- Onur Mutlu,

"Memory-Centric Computing Systems"

Invited Tutorial at *66th International Electron Devices Meeting (IEDM)*, Virtual, 12 December 2020.

[Slides (pptx) (pdf)]

[Executive Summary Slides (pptx) (pdf)]

[Tutorial Video (1 hour 51 minutes)]

[Executive Summary Video (2 minutes)]

[Abstract and Bio]

[Related Keynote Paper from VLSI-DAT 2020]

[Related Review Paper on Processing in Memory]

<https://www.youtube.com/watch?v=H3sEaINPBOE>



Memory-Centric Computing Systems

Onur Mutlu

omutlu@gmail.com

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

12 December 2020

IEDM Tutorial



SAFARI

ETH zürich

Carnegie Mellon



0:06 / 1:51:05



IEDM 2020 Tutorial: Memory-Centric Computing Systems, Onur Mutlu, 12 December 2020

1,641 views • Dec 23, 2020

48

0

SHARE

SAVE

...



Onur Mutlu Lectures
13.9K subscribers

<https://www.youtube.com/watch?v=H3sEaINPBOE>

ANALYTICS

EDIT VIDEO

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



Onur Mutlu



IEEE Custom Integrated Circuits Conference

Memory-Centric Computing

Onur Mutlu

Professor

ETH Zurich & Carnegie Mellon University

23 April 2023



6:27 / 1:35:22
🕒 ⏴ ⏵ ⏷ ⏸ ⏹ ⏺ ⏻ ⏻



IEEE CICC Educational Session Talk on Memory-Centric Computing (Prof. Onur Mutlu)



Onur Mutlu Lectures
36.6K subscribers

Subscribe

Like 24

Share

Clip

Save

...

526 views Streamed live on 23 Apr 2023

IEEE CICC Educational Session Talk on Memory-Centric Computing
Presenter: Professor Onur Mutlu (<https://people.inf.ethz.ch/omutlu/>)
Date: April 23, 2023

<https://www.youtube.com/watch?v=4x9nujJtqjM>

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



Seminar in Computer Architecture

Memory-Centric Computing

Geraldo F. Oliveira
Professor Onur Mutlu
ETH Zürich
Fall 2023
19 October 2023

SAFARI

1:45 / 1:54:48



Seminar in Computer Architecture - Lecture 4: Memory-Centric Computing (Fall 2023)



Onur Mutlu Lectures
36.6K subscribers

Subscribe

21



Share

Clip

Save



500 views Streamed live on 19 Oct 2023 [Livestream - Seminar in Computer Architecture - ETH Zürich \(Fall 2023\)](#)

Seminar in Computer Architecture, ETH Zürich, Fall 2023 (https://safari.ethz.ch/architecture_s...)

Lecture 4: Memory-Centric Computing

Lecturer: Prof. Onur Mutlu (<https://people.inf.ethz.ch/omutlu/>)

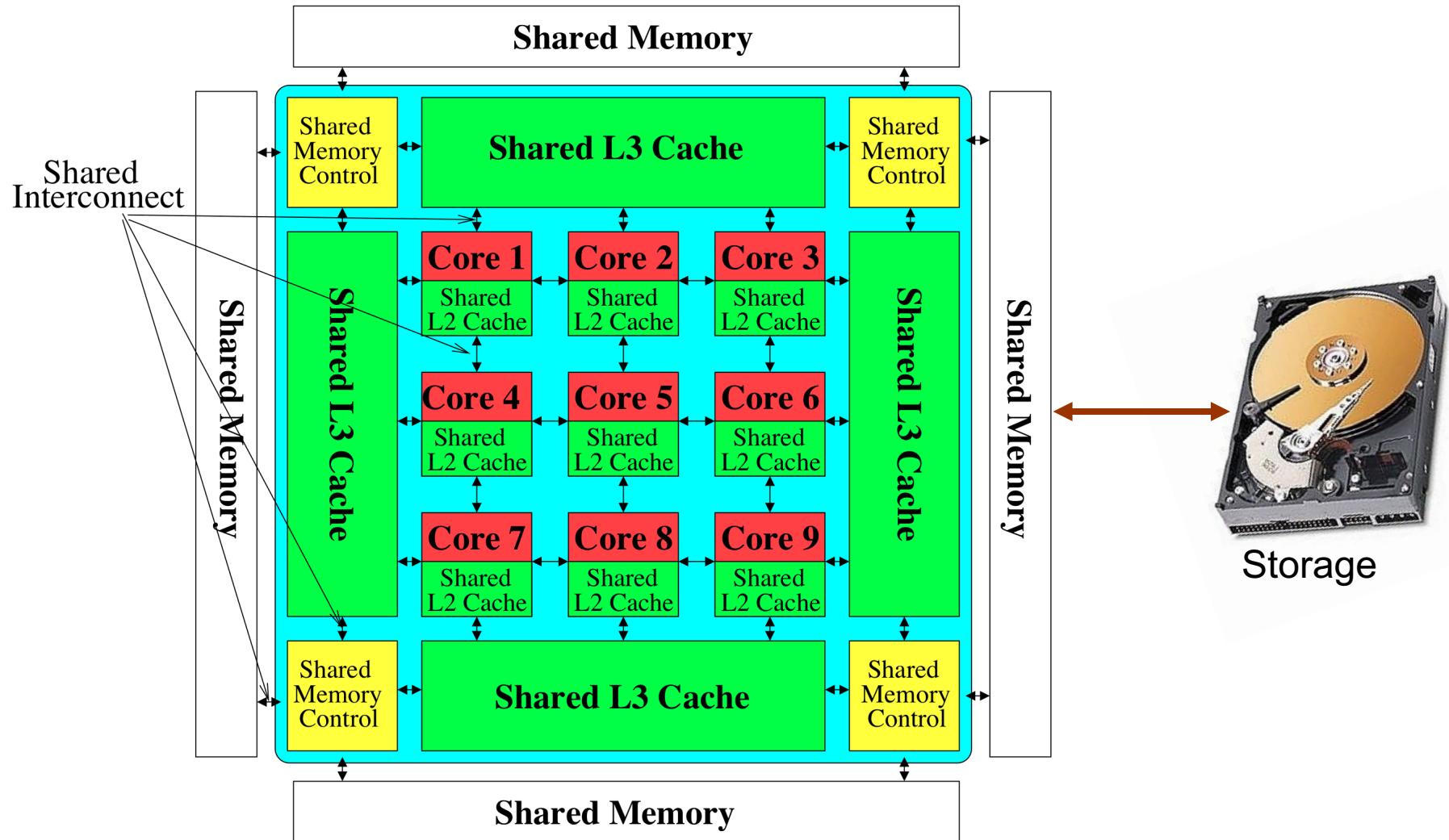
Date: October 19, 2023

<https://www.youtube.com/watch?v=CW4HxIGfofo>

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

Shared Resource Design for Multi-Core Systems

Memory System: A *Shared Resource* View



Most of the system is dedicated to storing and moving data

Resource Sharing Concept

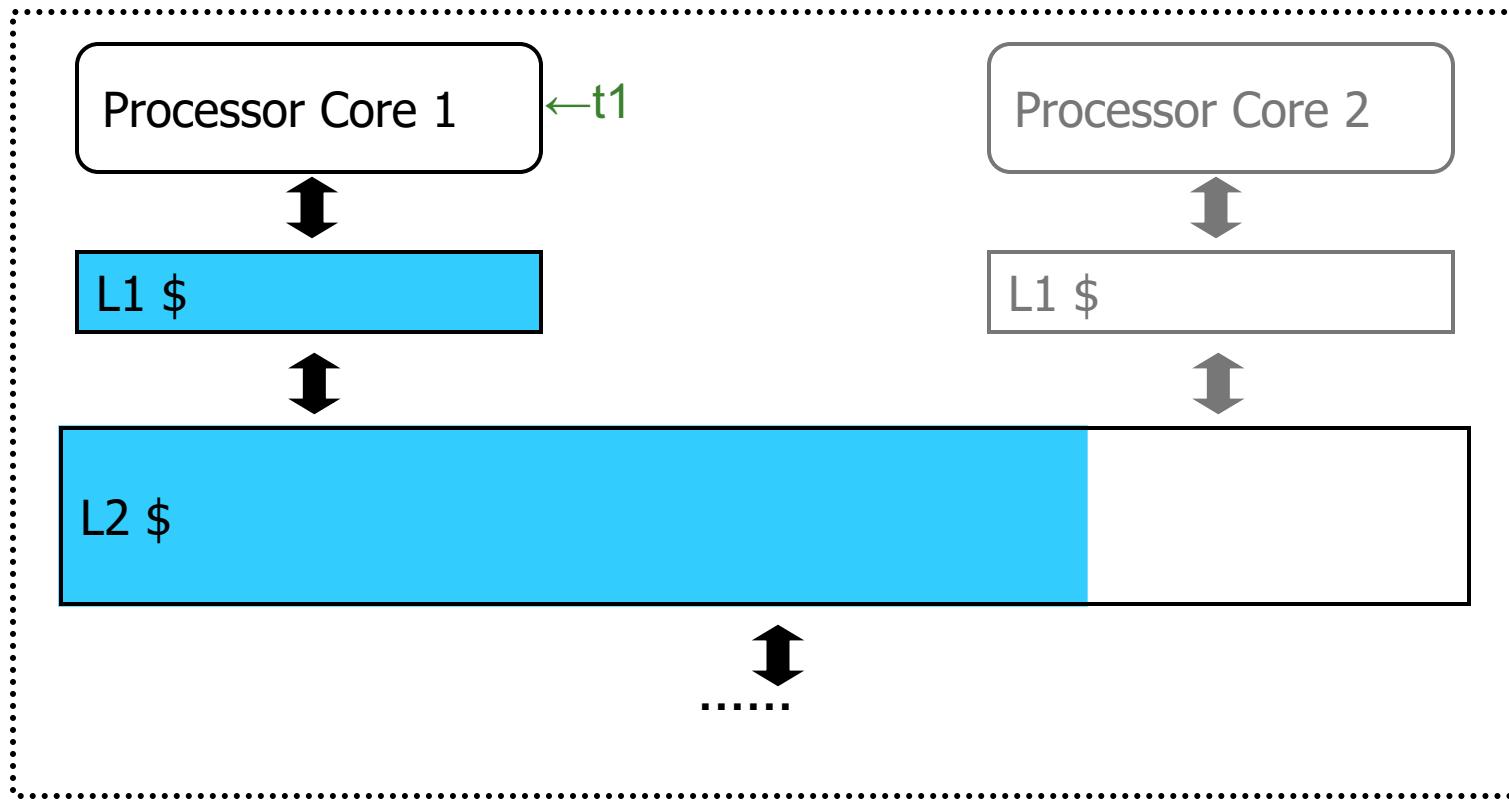
- Idea: Instead of dedicating a hardware resource to a hardware context, allow multiple contexts to use it
 - Example resources: functional units, pipeline, caches, buses, memory, interconnects, storage
- Why?
 - + Resource sharing improves utilization/efficiency → throughput
 - When a resource is left idle by one thread, another thread can use it; no need to replicate shared data
 - + Reduces communication latency
 - For example, shared data kept in the same cache in SMT processors
 - + Compatible with the shared memory model

Resource Sharing Disadvantages

- Resource sharing results in **contention for resources**
 - When the resource is not idle, another thread cannot use it
 - If space is occupied by one thread, another thread needs to re-occupy it
- Sometimes reduces each or some thread's performance
 - Thread performance can be worse than when it is run alone
- Eliminates performance isolation → inconsistent performance across runs
 - Thread performance depends on co-executing threads
- Uncontrolled (free-for-all) sharing degrades QoS
 - Causes unfairness, starvation

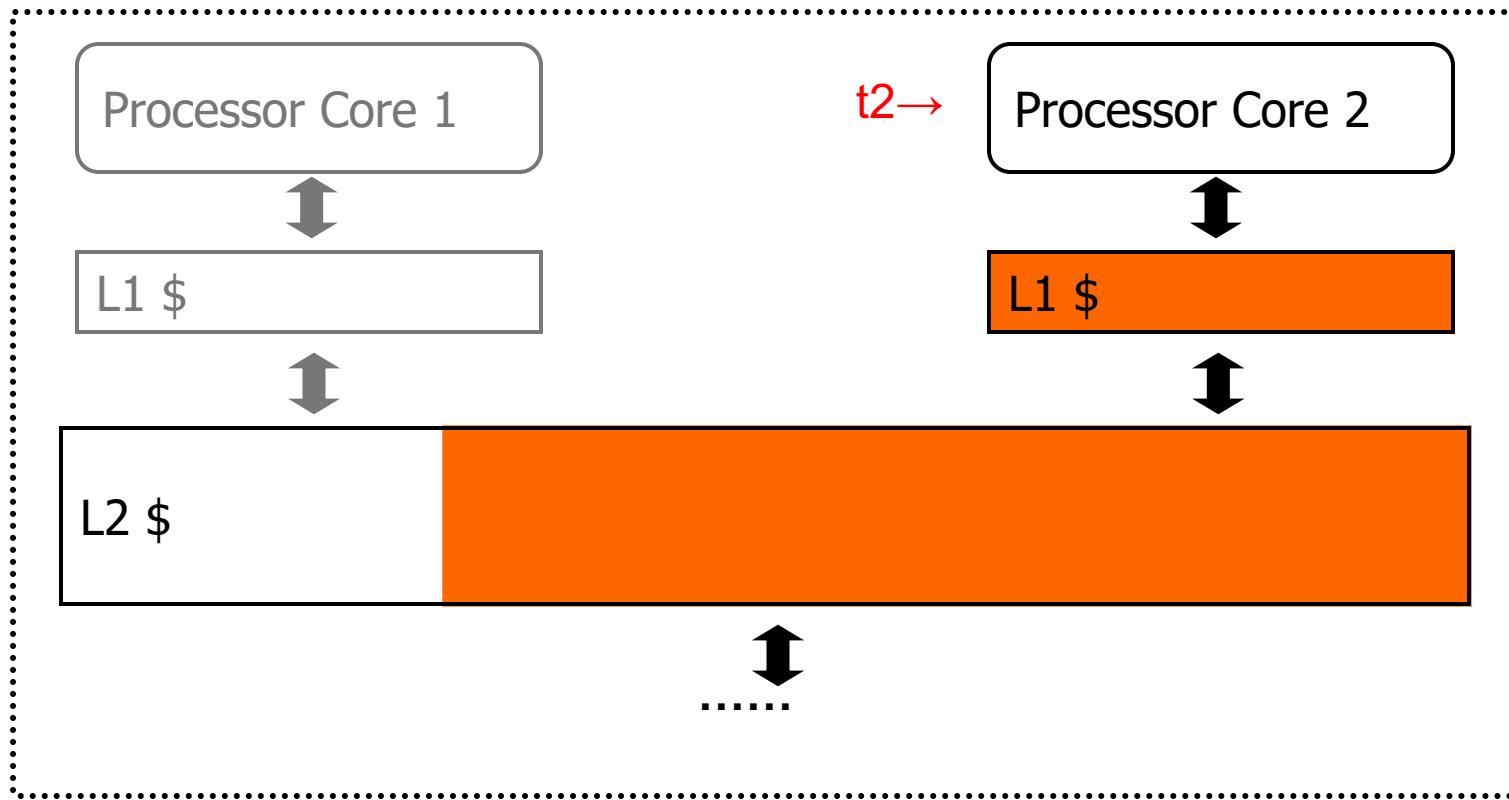
Need to efficiently and fairly utilize shared resources

Example: Problem with Shared Caches



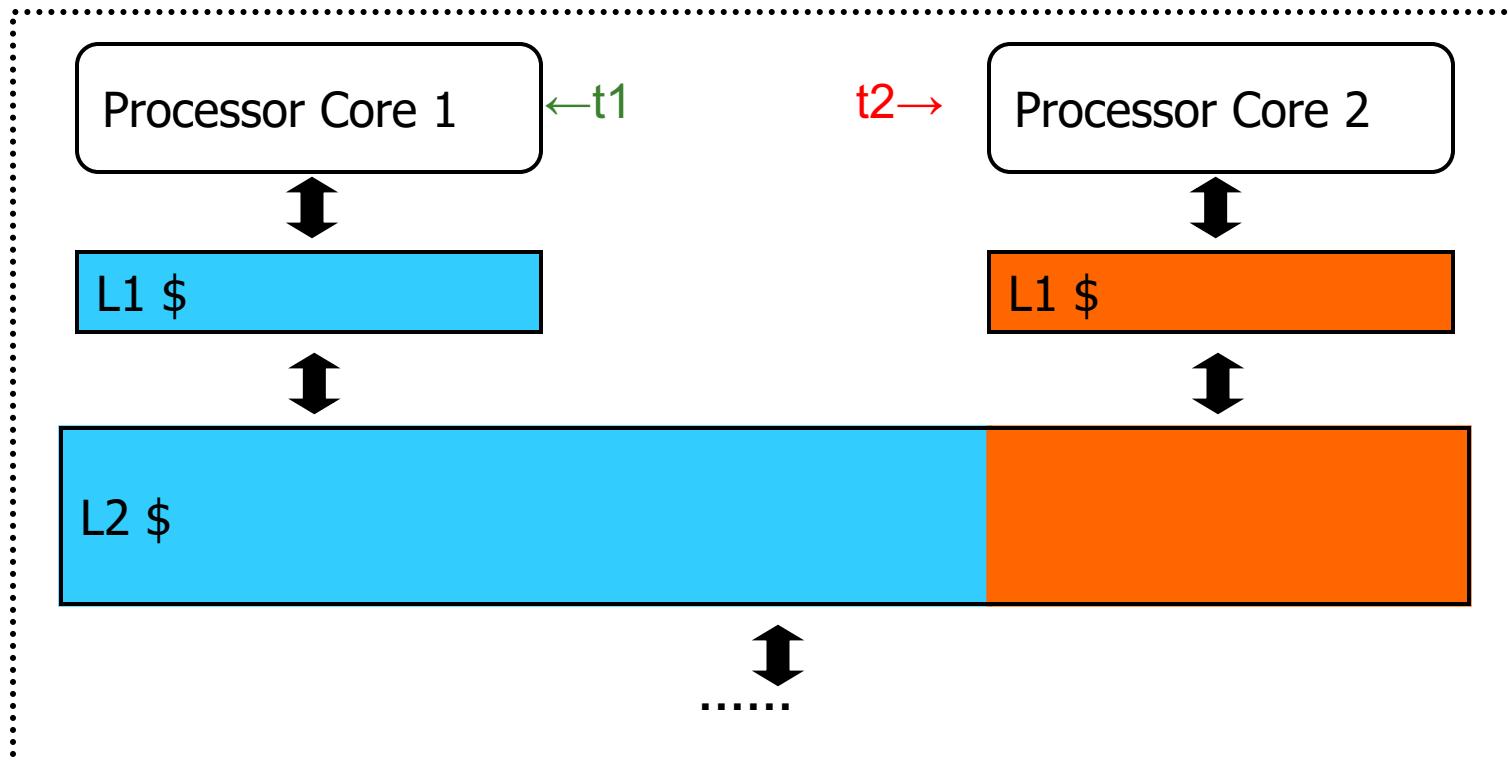
Kim et al., “Fair Cache Sharing and Partitioning in a Chip Multiprocessor Architecture,” PACT 2004.

Example: Problem with Shared Caches



Kim et al., “Fair Cache Sharing and Partitioning in a Chip Multiprocessor Architecture,” PACT 2004.

Example: Problem with Shared Caches



t2's throughput is significantly reduced due to unfair cache sharing.

Kim et al., “Fair Cache Sharing and Partitioning in a Chip Multiprocessor Architecture,” PACT 2004.

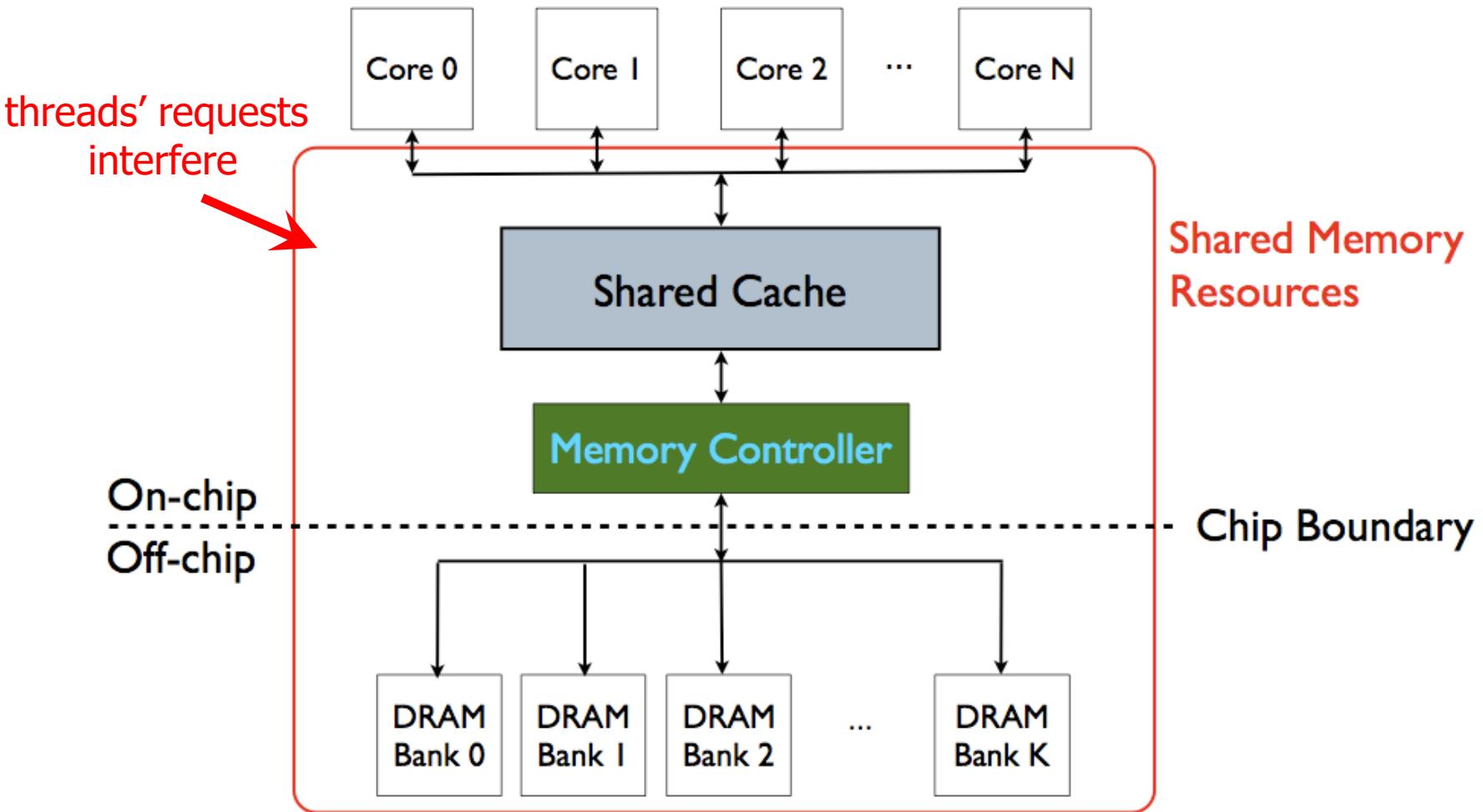
Need for QoS and Shared Resource Mgmt.

- Why is unpredictable performance (or lack of QoS) bad?
- Makes programmer's life difficult
 - An optimized program can get low performance (and performance varies widely depending on co-runners)
- Causes discomfort to user
 - An important program can starve
 - Examples from shared software resources
- Makes system management difficult
 - How do we enforce a Service Level Agreement when hardware resource sharing is uncontrollable?

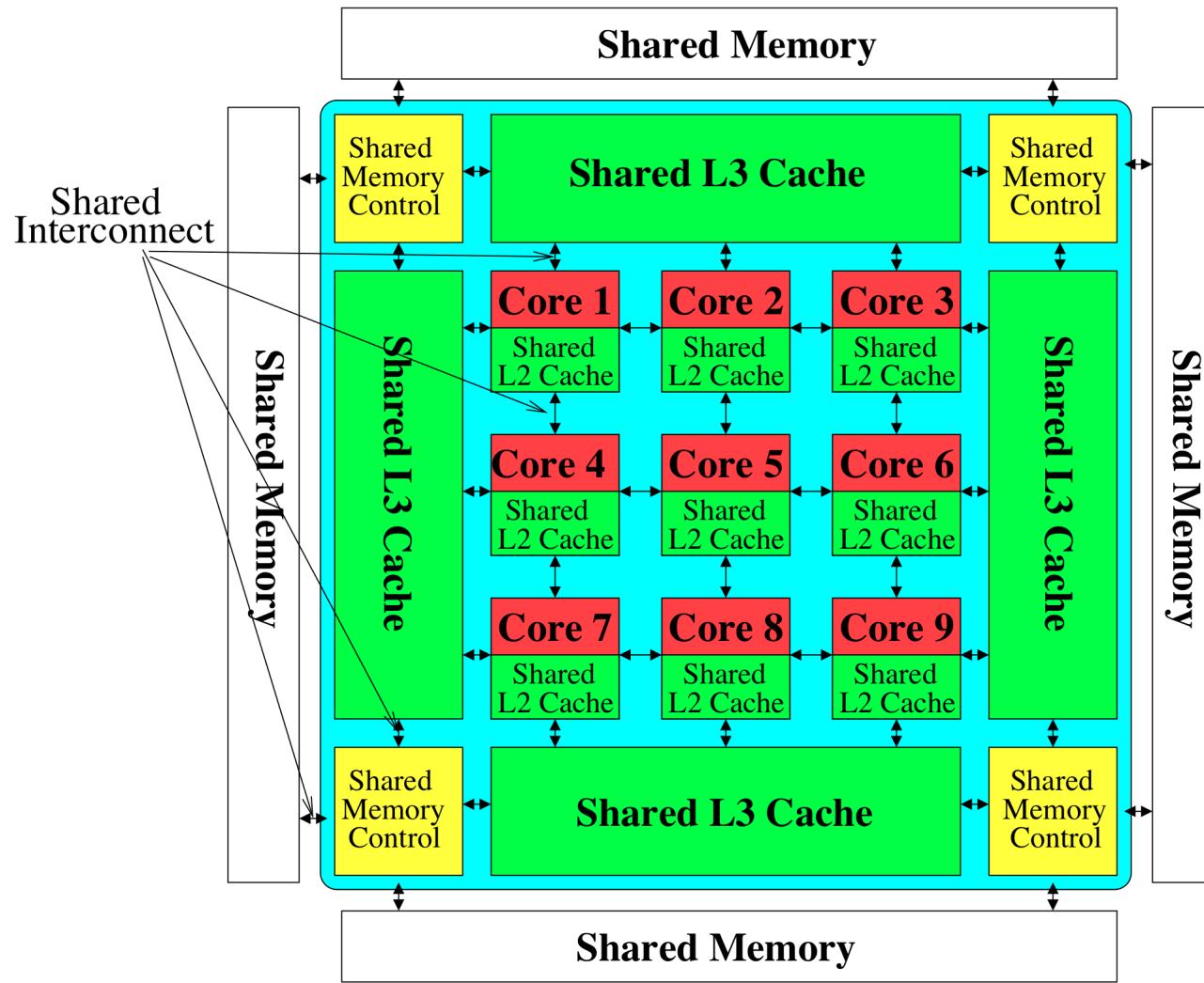
Resource Sharing vs. Partitioning

- Sharing improves throughput
 - Better utilization of space
- Partitioning provides performance isolation (predictable performance)
 - Dedicated space
- Can we get the benefits of both?
- Idea: Design shared resources such that they are efficiently utilized, controllable and partitionable
 - No wasted resource + QoS mechanisms for threads

Memory System is the Major Shared Resource



Much More of a Shared Resource in Future



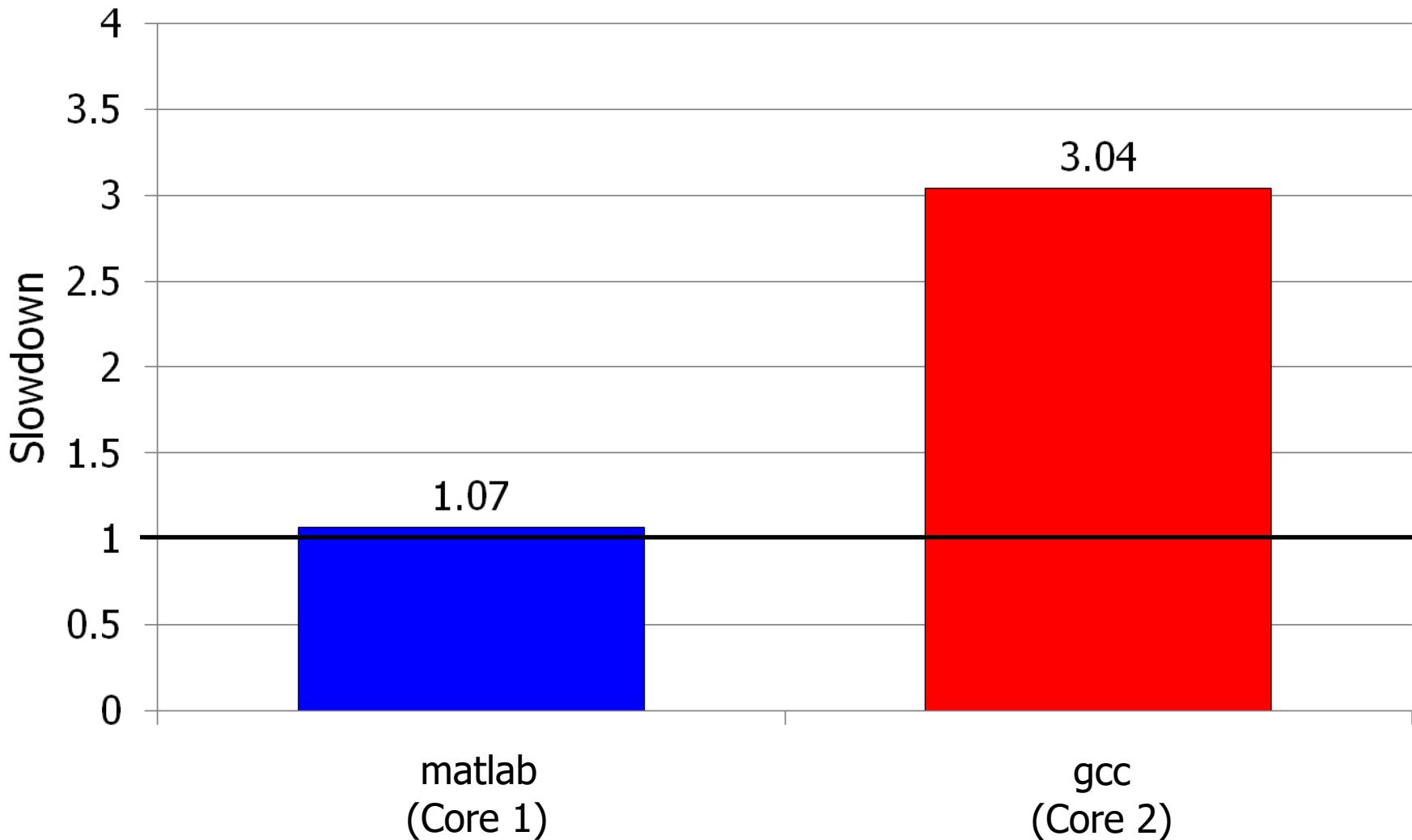
Most of the system is dedicated to storing and moving data

Inter-Thread/Application Interference

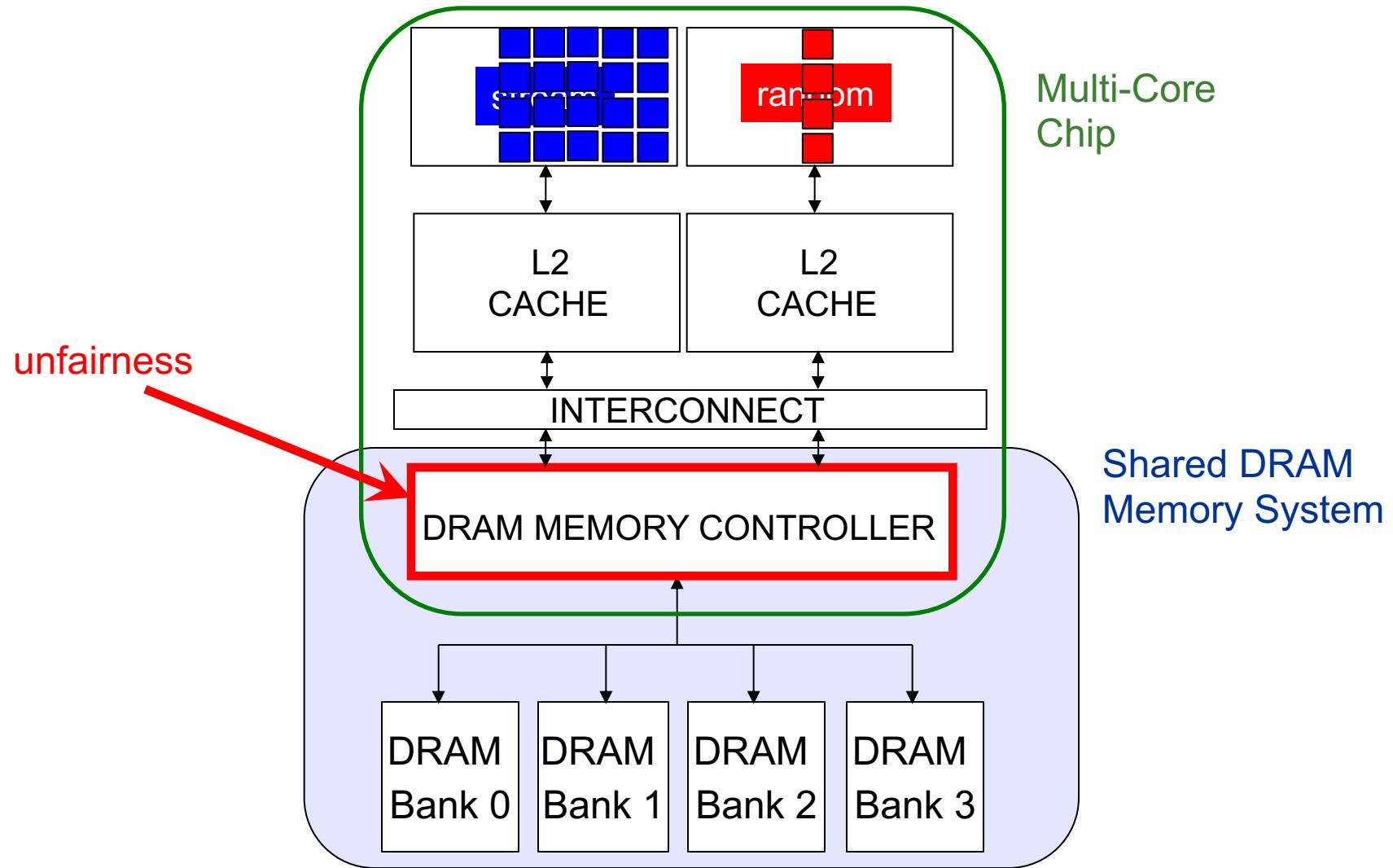
- Problem: Threads share the memory system, but memory system does not distinguish between threads' requests

- Existing memory systems
 - Free-for-all, shared based on demand
 - Control algorithms thread-unaware and thread-unfair
 - Aggressive threads can deny service to others
 - Do not try to reduce or control inter-thread interference

Unfair Slowdowns due to Interference



Uncontrolled Interference: An Example



A Memory Performance Hog

```
// initialize large arrays A, B  
  
for (j=0; j<N; j++) {  
    index = j*linesize; streaming  
    A[index] = B[index];  
    ...  
}
```

STREAM

- Sequential memory access
- Very high row buffer locality (96% hit rate)
- Memory intensive

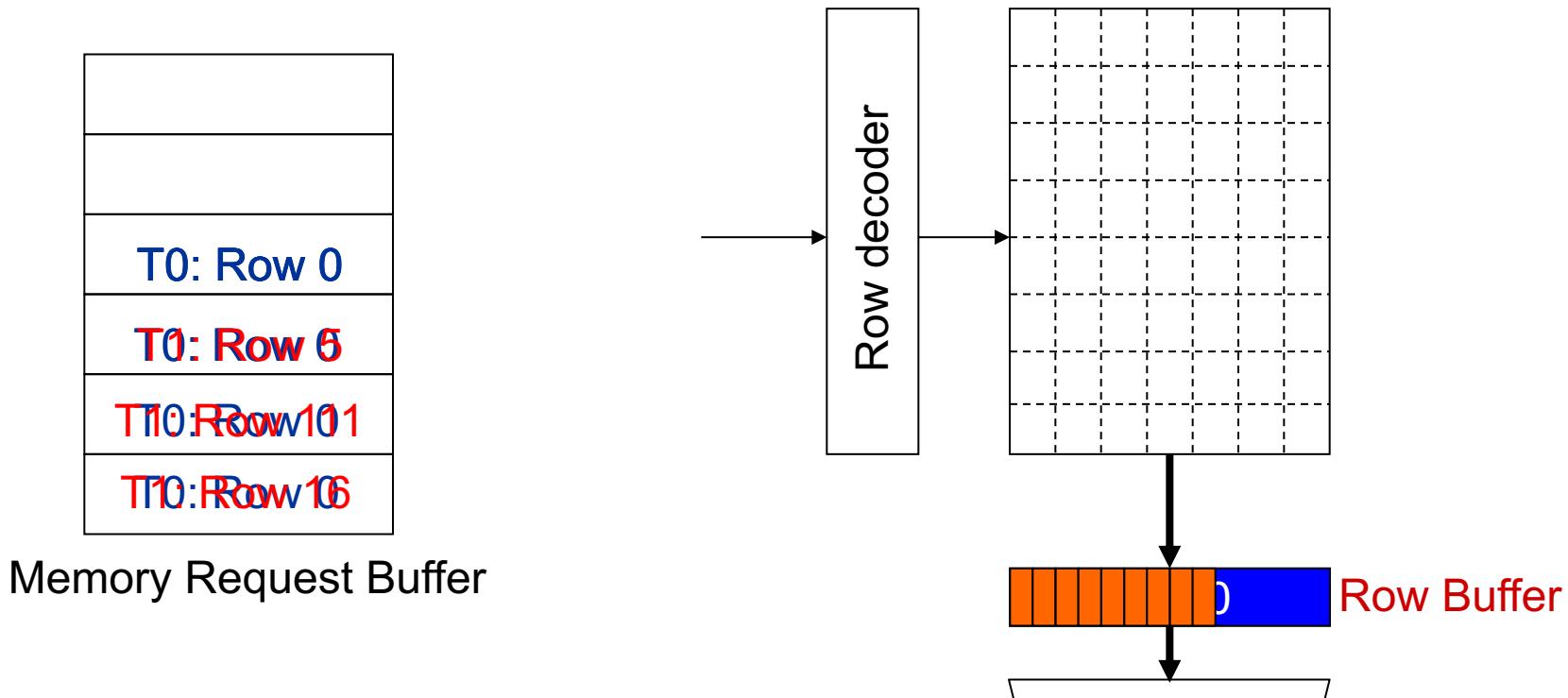
```
// initialize large arrays A, B  
  
for (j=0; j<N; j++) {  
    index = rand(); random  
    A[index] = B[index];  
    ...  
}
```

RANDOM

- Random memory access
- Very low row buffer locality (3% hit rate)
- Similarly memory intensive

Moscibroda and Mutlu, “Memory Performance Attacks,” USENIX Security 2007.

What Does the Memory Hog Do?



Row size: 8KB, cache block size: 64B

128 (8KB/64B) requests of T0 serviced before T1

Moscibroda and Mutlu, “Memory Performance Attacks,” USENIX Security 2007.

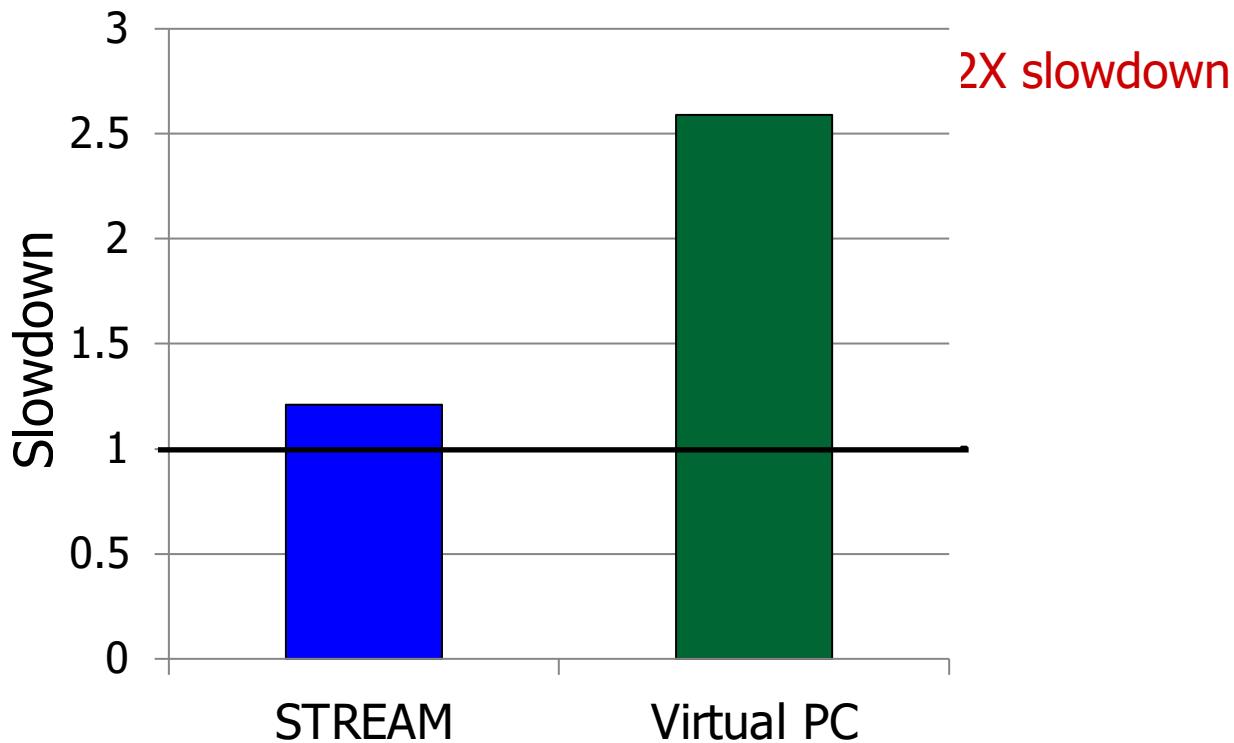
DRAM Controllers

- A row-conflict memory access takes significantly longer than a row-hit access
- Current controllers take advantage of the row buffer
- Commonly used scheduling policy (FR-FCFS) [Rixner 2000]*
 - (1) Row-hit first: Service row-hit memory accesses first
 - (2) Oldest-first: Then service older accesses first
- This scheduling policy aims to maximize DRAM throughput
 - But, it is unfair when multiple threads share the DRAM system

*Rixner et al., “Memory Access Scheduling,” ISCA 2000.

*Zuravleff and Robinson, “Controller for a synchronous DRAM ...,” US Patent 5,630,096, May 1997.

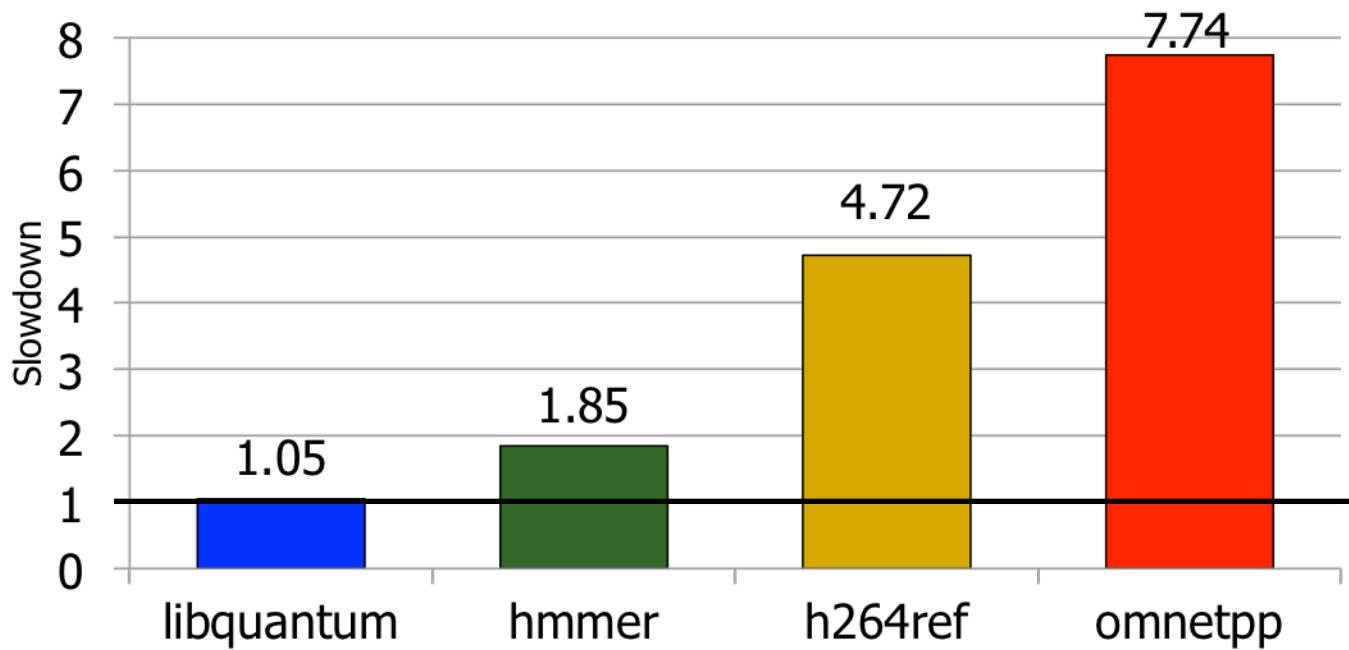
Effect of the Memory Performance Hog



Results on Intel Pentium D running Windows XP
(Similar results for Intel Core Duo and AMD Turion, and on Fedora Linux)

Moscibroda and Mutlu, “Memory Performance Attacks,” USENIX Security 2007.

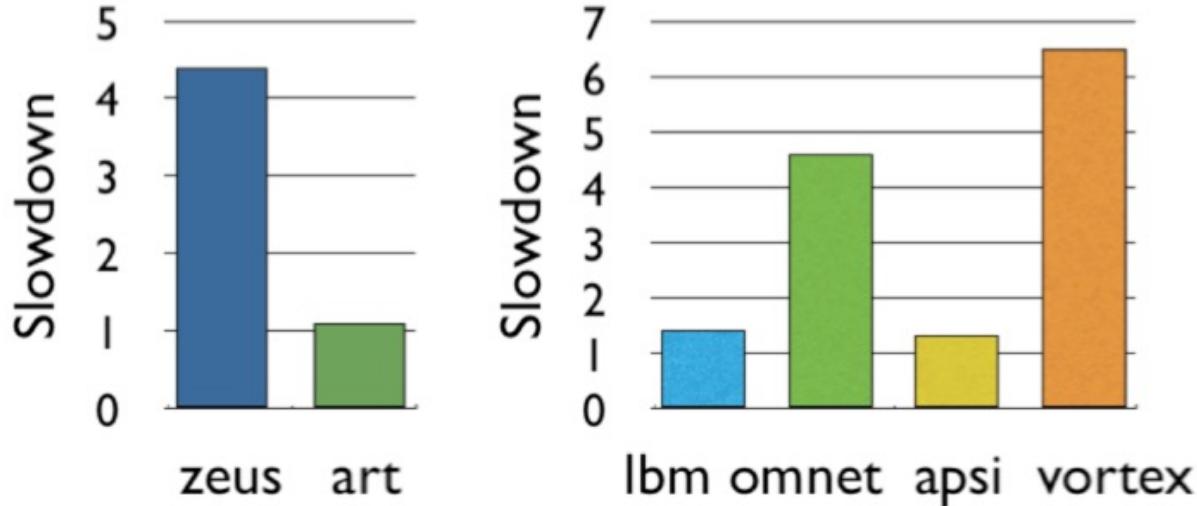
Greater Problem with More Cores



- Vulnerable to denial of service (DoS)
- Unable to enforce priorities or SLAs
- Low system performance

Uncontrollable, unpredictable system

Greater Problem with More Cores



- Vulnerable to denial of service (DoS)
- Unable to enforce priorities or SLAs
- Low system performance

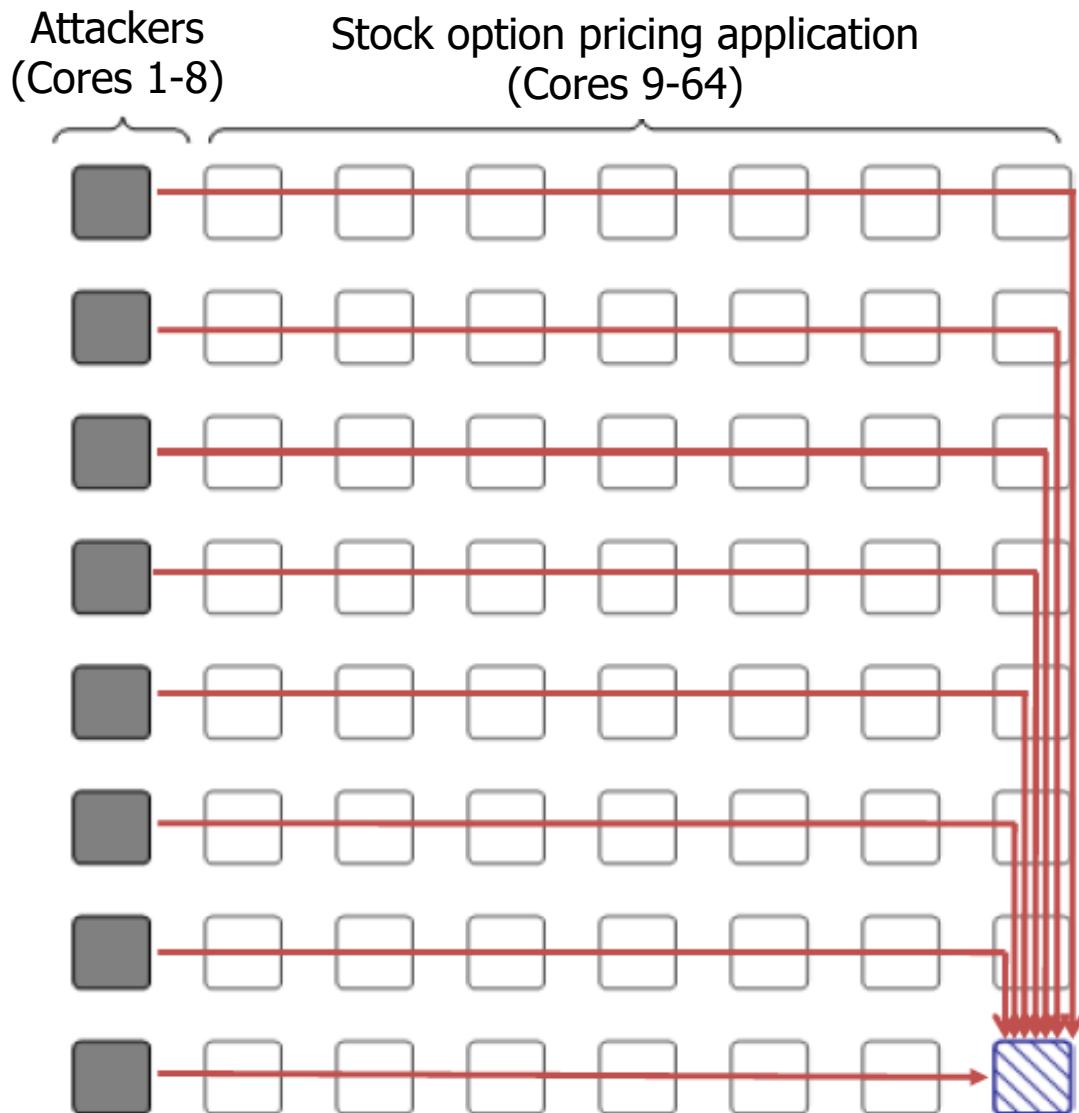
Uncontrollable, unpredictable system

Distributed DoS in Networked Multi-Core Systems

Cores connected via
packet-switched
routers on chip

~5000X latency increase

Grot, Hestness, Keckler, Mutlu,
"Preemptive virtual clock: A Flexible,
Efficient, and Cost-effective QOS
Scheme for Networks-on-Chip,"
MICRO 2009.



More on Memory Performance Attacks

- 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

More on Interconnect Based Starvation

- Boris Grot, Stephen W. Keckler, and Onur Mutlu,
"Preemptive Virtual Clock: A Flexible, Efficient, and Cost-effective QOS Scheme for Networks-on-Chip"
Proceedings of the 42nd International Symposium on Microarchitecture (MICRO), pages 268-279, New York, NY, December 2009. [Slides \(pdf\)](#)

Preemptive Virtual Clock: A Flexible, Efficient, and Cost-effective QOS Scheme for Networks-on-Chip

Boris Grot

Stephen W. Keckler

Onur Mutlu[†]

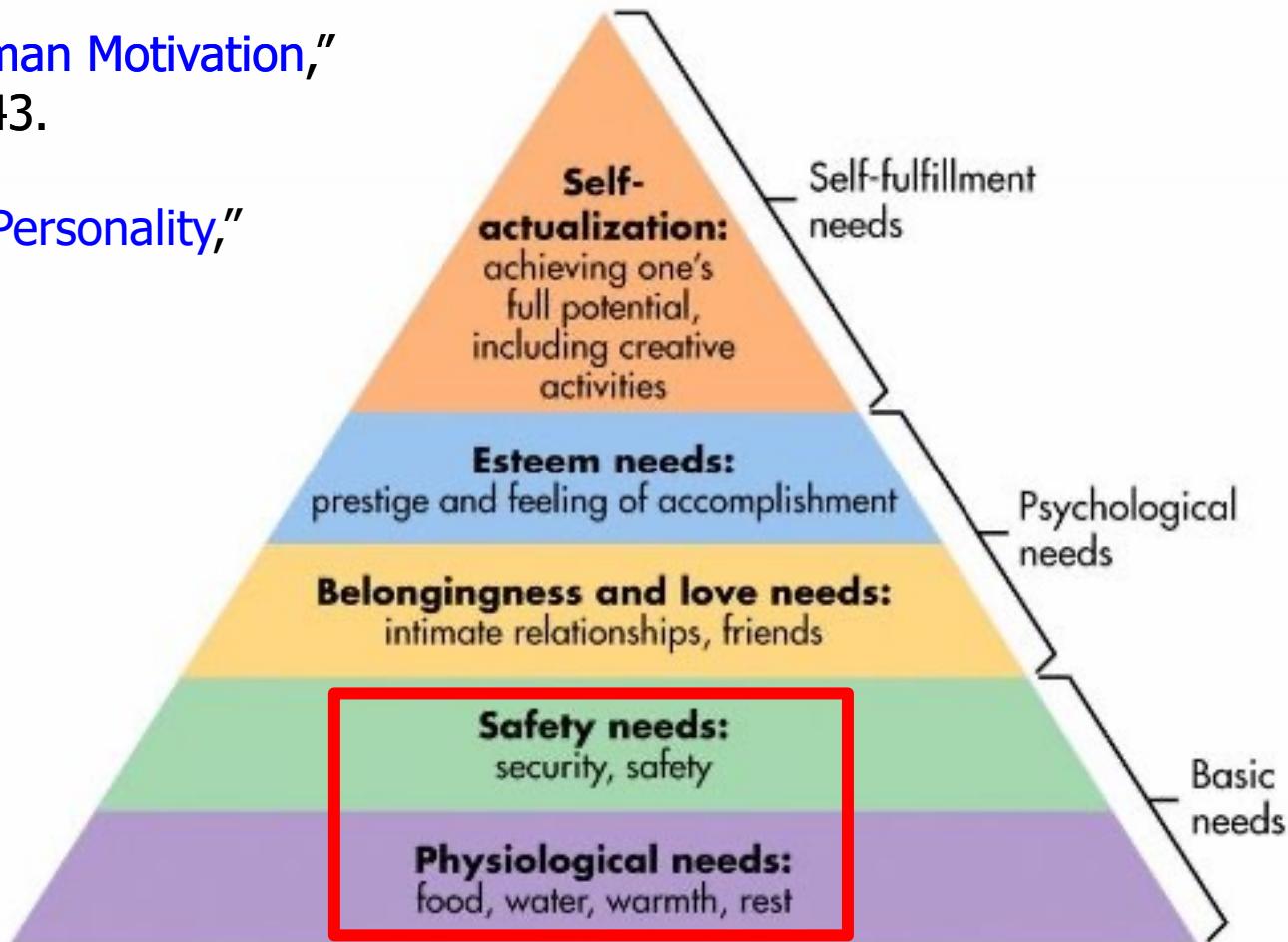
Department of Computer Sciences
The University of Texas at Austin
{bgrot, skeckler}@cs.utexas.edu}

[†]Computer Architecture Laboratory (CALCM)
Carnegie Mellon University
onur@cmu.edu

Maslow's (Human) Hierarchy of Needs

Maslow, "A Theory of Human Motivation,"
Psychological Review, 1943.

Maslow, "Motivation and Personality,"
Book, 1954-1970.



- Lack of QoS can be a **safety and security** problem

How Do We Solve The Problem?

- Inter-thread interference is uncontrolled in all memory resources
 - Memory controller
 - Interconnect
 - Caches
- We need to control it
 - i.e., design an interference-aware (QoS-aware) memory system

QoS-Aware Memory Systems: Challenges

- How do we **reduce inter-thread interference?**
 - Improve system performance and core utilization
 - Reduce request serialization and core starvation
- How do we **control inter-thread interference?**
 - Provide mechanisms to enable system software to enforce QoS policies
 - While providing high system performance
- How do we **make the memory system configurable/flexible?**
 - Enable flexible mechanisms that can achieve many goals
 - Provide fairness or throughput when needed
 - Satisfy performance guarantees when needed

Designing QoS-Aware Memory Systems: Approaches

- **Smart resources:** Design each shared resource to have a configurable interference control/reduction mechanism
 - QoS-aware memory controllers
 - QoS-aware interconnects
 - QoS-aware caches
- **Dumb resources:** Keep each resource free-for-all, but reduce/control interference by injection control or data mapping
 - Source throttling to control access to memory system
 - QoS-aware data mapping to memory controllers
 - QoS-aware thread scheduling to cores

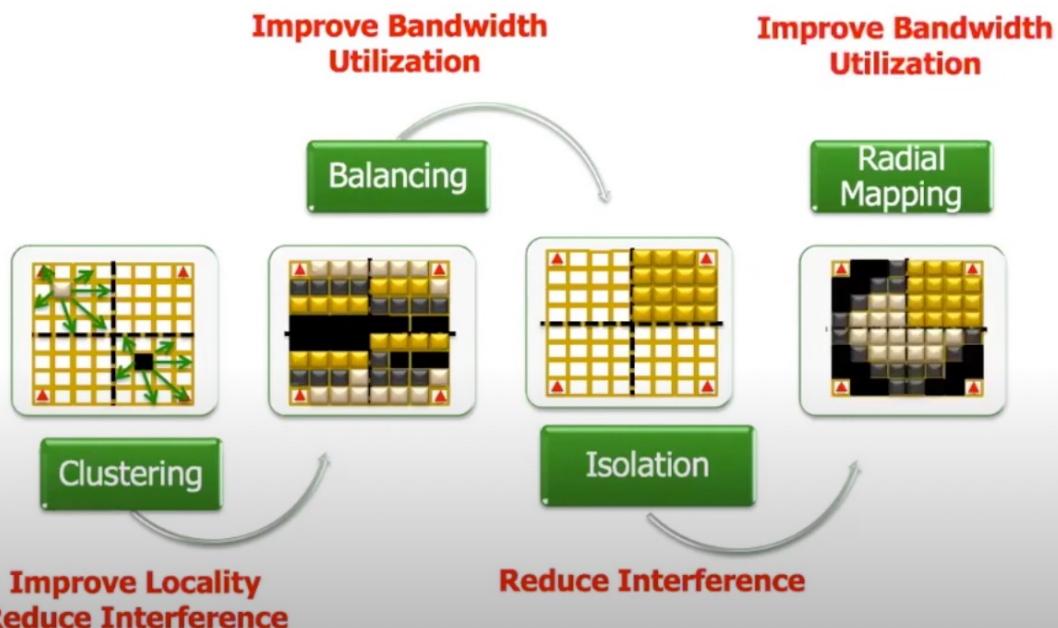
Fundamental Interference Control Techniques

- Goal: to reduce/control inter-thread memory interference

1. Prioritization or request scheduling
2. Data mapping to banks/channels/ranks
3. Core/source throttling
4. Application/thread scheduling

Lecture on Other QoS Techniques

Application-to-Core Mapping



ETH ZENTRUM

Computer Arch - Lecture 13: Memory Interference and Quality of Service II (ETH Zürich, Spring 2020)

1,006 views • Nov 7, 2020

89

CC G S D E

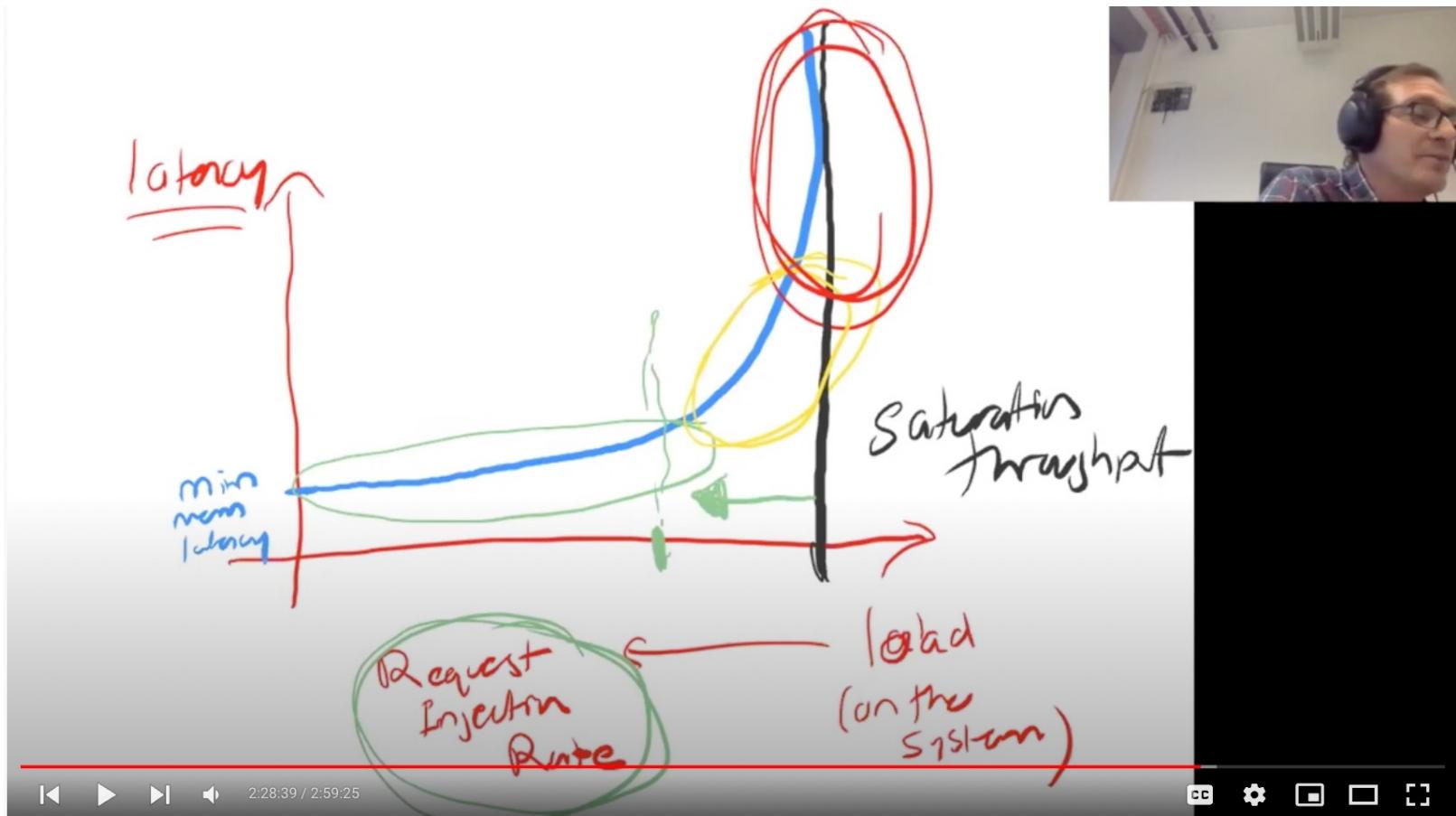


Onur Mutlu Lectures
19.8K subscribers

ANALYTICS

EDIT VIDEO

Lecture on Other QoS Techniques



ETH ZENTRUM

Computer Arch - Lecture 13: Memory Interference and Quality of Service II (ETH Zürich, Spring 2020)

1,006 views • Nov 7, 2020

23 likes 0 dislikes SHARE SAVE ...



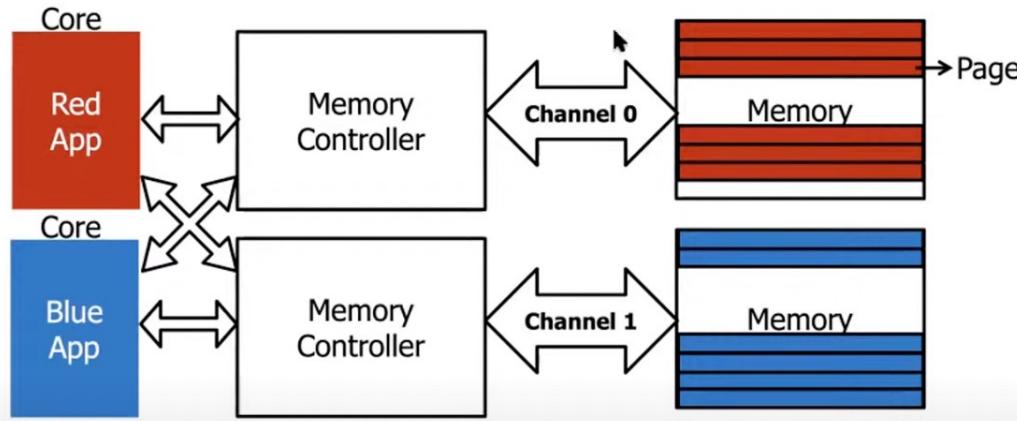
Onur Mutlu Lectures
19.8K subscribers

ANALYTICS

EDIT VIDEO

Memory Channel Partitioning

Partitioning Channels Between Applications



Eliminates interference between applications' requests

◀ ▶ ⏪ 21:30 / 1:20:55 ⏩

zoom
CC 29 HD

ETH ZURICH D-ITET

Seminar in Computer Architecture - Lecture 4: Memory Channel Partitioning (Fall 2021)

379 views • Streamed live on Oct 14, 2021

19 0 SHARE SAVE ...



Onur Mutlu Lectures
19.8K subscribers

ANALYTICS

EDIT VIDEO

Handling CPU-IO Interference

Donghyuk Lee, Lavanya Subramanian, Rachata Ausavarungnirun, Jongmoo Choi,
and Onur Mutlu,

"Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM"

Proceedings of the 24th International Conference on Parallel Architectures and Compilation Techniques (PACT), San Francisco, CA, USA, October 2015.

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

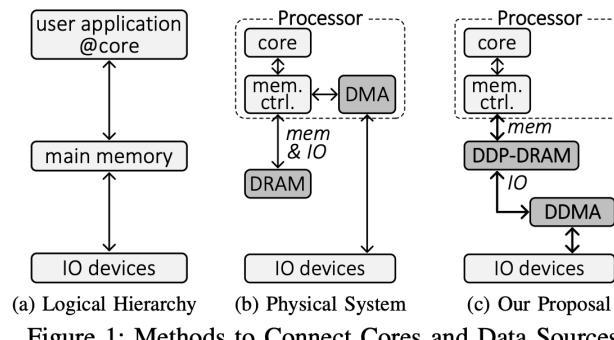


Figure 1: Methods to Connect Cores and Data Sources

Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM

Donghyuk Lee* Lavanya Subramanian* Rachata Ausavarungnirun* Jongmoo Choi† Onur Mutlu*

*Carnegie Mellon University

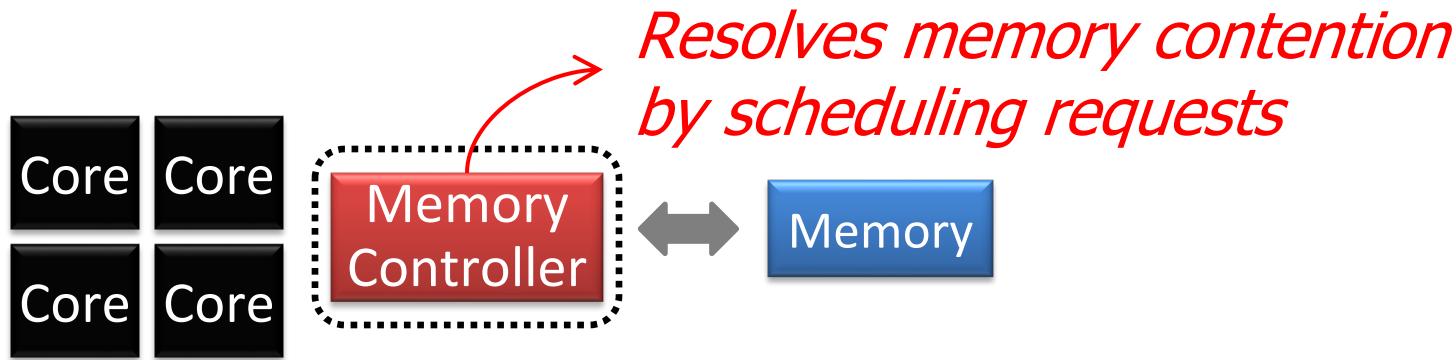
{donghyu1, lsubrama, rachata, onur}@cmu.edu

†Dankook University

choijm@dankook.ac.kr

QoS-Aware Memory Scheduling: Revolution & Evolution

QoS-Aware Memory Scheduling



- How to schedule requests to provide
 - High system performance
 - High fairness to applications
 - Configurability to system software
- Memory controller needs to be aware of threads

QoS-Aware Memory Scheduling: Evolution

- **Stall-time fair memory scheduling** [Mutlu+ MICRO'07]
 - Idea: Estimate and balance thread slowdowns
 - Takeaway: **Proportional thread progress improves performance, especially when threads are “heavy”** (memory intensive)
- **Parallelism-aware batch scheduling** [Mutlu+ ISCA'08, Top Picks'09]
 - Idea: Rank threads and service in rank order (to preserve bank parallelism); batch requests to prevent starvation
 - Takeaway: **Preserving within-thread bank-parallelism improves performance**; request batching improves fairness
- **ATLAS memory scheduler** [Kim+ HPCA'10]
 - Idea: Prioritize threads that have attained the least service from the memory scheduler
 - Takeaway: **Prioritizing “light” threads improves performance**

QoS-Aware Memory Scheduling: Evolution

- Thread cluster memory scheduling [Kim+ MICRO'10, Top Picks'11]
 - Idea: Cluster threads into two groups (latency vs. bandwidth sensitive); prioritize the latency-sensitive ones; employ a fairness policy in the bandwidth sensitive group
 - Takeaway: **Heterogeneous scheduling policy that is different based on thread behavior maximizes both performance and fairness**
- Integrated Memory Channel Partitioning and Scheduling [Muralidhara+ MICRO'11]
 - Idea: Only prioritize very latency-sensitive threads in the scheduler; mitigate all other applications' interference via channel partitioning
 - Takeaway: **Intelligently combining application-aware channel partitioning and memory scheduling provides better performance than either**

QoS-Aware Memory Scheduling: Evolution

- Parallel application memory scheduling [Ebrahimi+ MICRO'11]
 - Idea: Identify and prioritize limiter threads of a multithreaded application in the memory scheduler; provide fast and fair progress to non-limiter threads
 - Takeaway: Carefully prioritizing between limiter and non-limiter threads of a parallel application improves performance
- Staged memory scheduling [Ausavarungnirun+ ISCA'12]
 - Idea: Divide the functional tasks of an application-aware memory scheduler into multiple distinct stages, where each stage is significantly simpler than a monolithic scheduler
 - Takeaway: Staging enables the design of a scalable and relatively simpler application-aware memory scheduler that works on very large request buffers

QoS-Aware Memory Scheduling: Evolution

- MISE: Memory Slowdown Model [Subramanian+ HPCA'13]
 - Idea: Estimate the performance of a thread by estimating its change in memory request service rate when run alone vs. shared → use this simple model to estimate slowdown to design a scheduling policy that provides predictable performance or fairness
 - Takeaway: Request service rate of a thread is a good proxy for its performance; alone request service rate can be estimated by giving high priority to the thread in memory scheduling for a while
- ASM: Application Slowdown Model [Subramanian+ MICRO'15]
 - Idea: Extend MISE to take into account cache+memory interference
 - Takeaway: Cache access rate of an application can be estimated accurately and is a good proxy for application performance

QoS-Aware Memory Scheduling: Evolution

- **BLISS: Blacklisting Memory Scheduler** [Subramanian+ ICCD'14, TPDS'16]
 - Idea: Deprioritize (i.e., blacklist) a thread that has consecutively serviced a large number of requests
 - Takeaway: **Blacklisting greatly reduces interference enables the scheduler to be simple without requiring full thread ranking**
- **DASH: Deadline-Aware Memory Scheduler** [Usui+ TACO'16]
 - Idea: Balance prioritization between CPUs, GPUs and Hardware Accelerators (HWA) by keeping HWA progress in check vs. deadlines such that HWAs do not hog performance and appropriately distinguishing between latency-sensitive vs. bandwidth-sensitive CPU workloads
 - Takeaway: **Proper control of HWA progress and application-aware CPU prioritization leads to better system performance while meeting HWA deadlines**

QoS-Aware Memory Scheduling: Evolution

- Prefetch-aware shared resource management [Ebrahimi+ ISCA'11] [Ebrahimi+ MICRO'09] [Ebrahimi+ HPCA'09] [Lee+ MICRO'08'09]
 - Idea: Prioritize prefetches depending on how they affect system performance; even accurate prefetches can degrade performance of the system
 - Takeaway: Carefully controlling and prioritizing prefetch requests improves performance and fairness
- DRAM-Aware last-level cache policies and write scheduling [Lee+ HPS Tech Report'10] [Seshadri+ ISCA'14]
 - Idea: Design cache eviction and replacement policies such that they proactively exploit the state of the memory controller and DRAM (e.g., proactively evict data from the cache that hit in open rows)
 - Takeaway: Coordination of last-level cache and DRAM policies improves performance and fairness; writes should not be ignored

QoS-Aware Memory Scheduling: Evolution

- **FIRM: Memory Scheduling for NVM** [Zhao+ MICRO'14]
 - Idea: Carefully handle write-read prioritization with coarse-grained batching and application-aware scheduling
 - Takeaway: Carefully controlling and prioritizing write requests improves performance and fairness; write requests are especially critical in NVMs
- **Criticality-Aware Memory Scheduling for GPUs** [Jog+ SIGMETRICS'16]
 - Idea: Prioritize latency-critical cores' requests in a GPU system
 - Takeaway: Need to carefully balance locality and criticality to make sure performance improves by taking advantage of both
- **Worst-case Execution Time Based Memory Scheduling for Real-Time Systems** [Kim+ RTAS'14, JRTS'16]

Stall-Time Fair Memory Scheduling

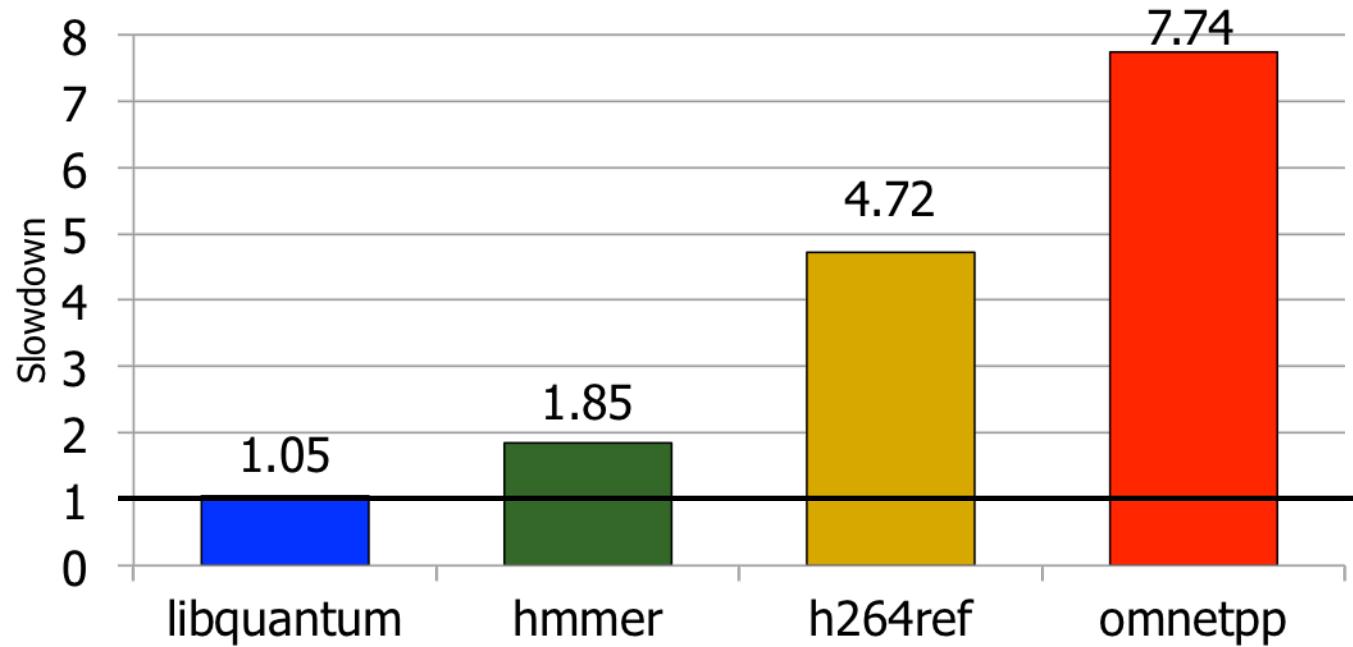
Onur Mutlu and Thomas Moscibroda,

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

40th International Symposium on Microarchitecture (MICRO),

pages 146-158, Chicago, IL, December 2007. Slides (ppt)

The Problem: Unfairness



- Vulnerable to denial of service (DoS)
- Unable to enforce priorities or SLAs
- Low system performance

Uncontrollable, unpredictable system

How Do We Solve the Problem?

- Stall-time fair memory scheduling [Mutlu+ MICRO'07]
- Goal: Threads sharing main memory should experience similar slowdowns compared to when they are run alone → fair scheduling
 - Also improves overall system performance by ensuring cores make “proportional” progress
- Idea: Memory controller estimates each thread’s slowdown due to interference and schedules requests in a way to balance the slowdowns
- Mutlu and Moscibroda, “[Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors](#),” MICRO 2007.

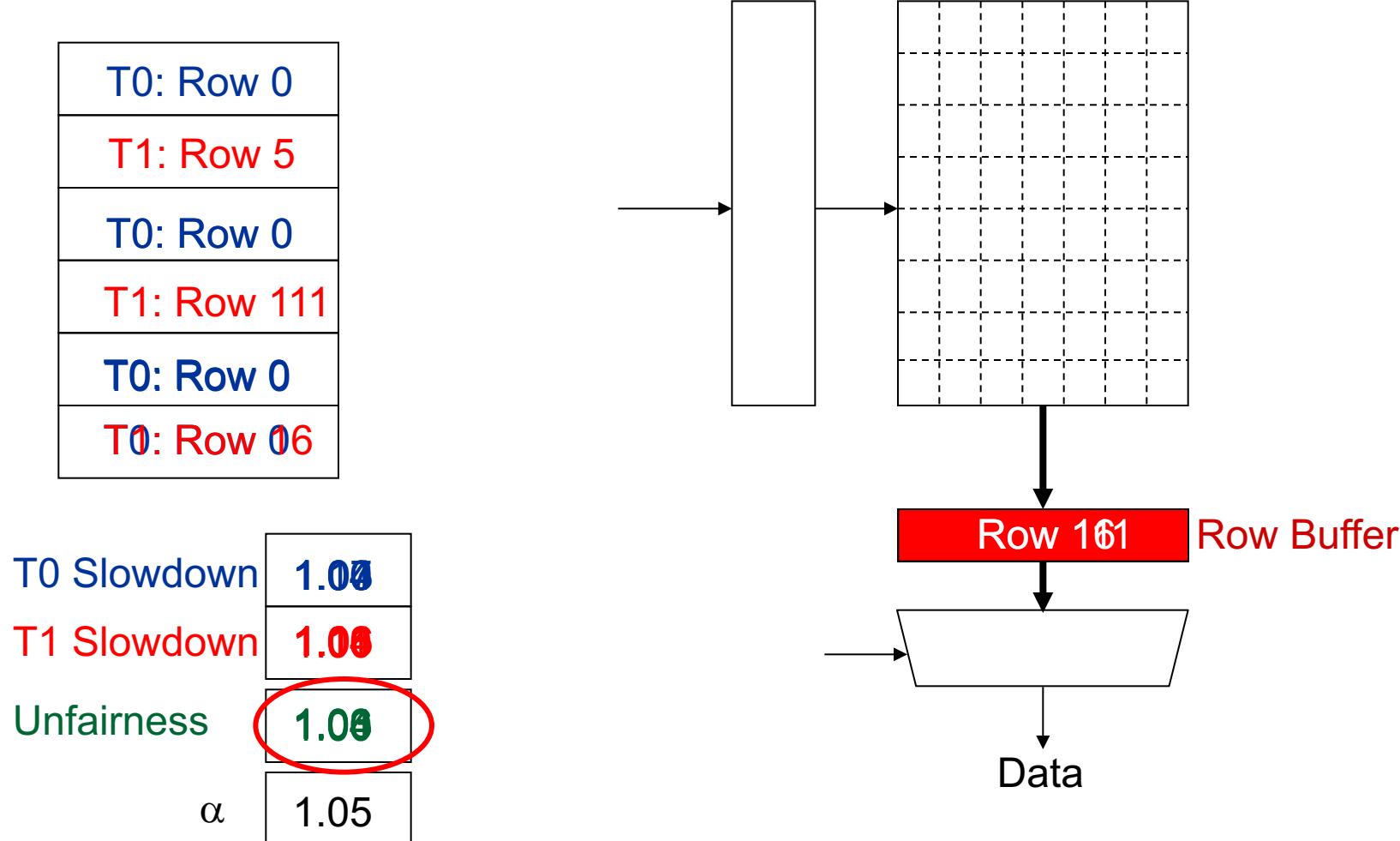
Stall-Time Fairness in Shared DRAM Systems

- A DRAM system is fair if it equalizes the slowdown of equal-priority threads relative to when each thread is run alone on the same system
- DRAM-related stall-time: The time a thread spends waiting for DRAM memory
- ST_{shared} : DRAM-related stall-time when the thread runs with other threads
- ST_{alone} : DRAM-related stall-time when the thread runs alone
- **Memory-slowdown = ST_{shared}/ST_{alone}**
 - Relative increase in stall-time
- *Stall-Time Fair Memory scheduler (STFM)* aims to **equalize** Memory-slowdown for interfering threads, without sacrificing performance
 - Considers inherent DRAM performance of each thread
 - Aims to allow proportional progress of threads

STFM Scheduling Algorithm [MICRO'07]

- For each thread, the DRAM controller
 - Tracks ST_{shared}
 - Estimates ST_{alone}
- Each cycle, the DRAM controller
 - Computes Slowdown = $ST_{\text{shared}}/ST_{\text{alone}}$ for threads with legal requests
 - Computes **unfairness** = MAX Slowdown / MIN Slowdown
- If unfairness < α
 - Use DRAM throughput oriented scheduling policy
- If **unfairness** $\geq \alpha$
 - Use fairness-oriented scheduling policy
 - (1) requests from thread with MAX Slowdown first
 - (2) row-hit first , (3) oldest-first

How Does STFM Prevent Unfairness?



STFM Pros and Cons

- Upsides:
 - First algorithm specialized for fair multi-core memory scheduling
 - Provides a mechanism to estimate memory slowdown of a thread
 - Good at providing fairness
 - Being fair can improve performance

- Downsides:
 - Does not handle all types of interference
 - Complex to implement
 - Slowdown estimations can be incorrect

More on STFM

- 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

Parallelism-Aware Batch Scheduling

Onur Mutlu and Thomas Moscibroda,

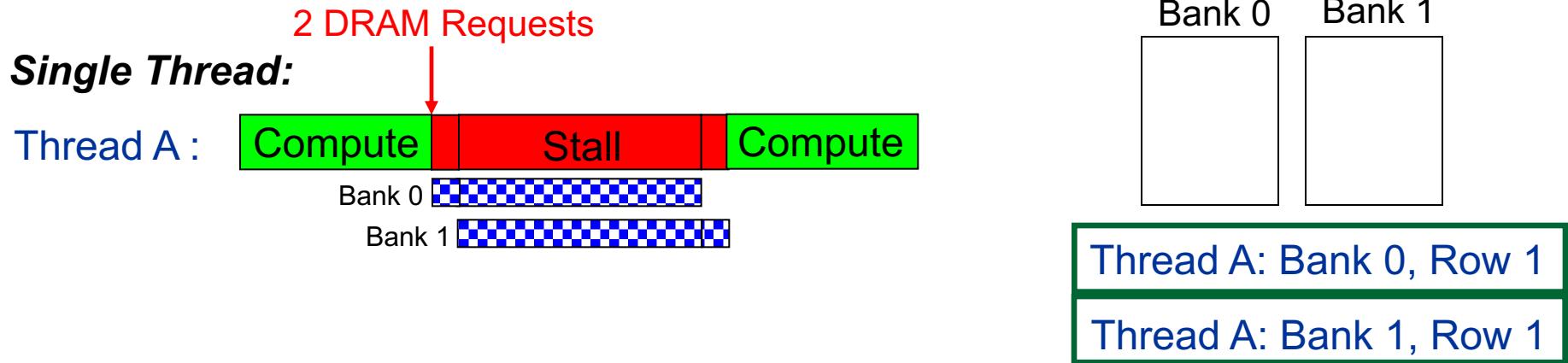
"Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems"

35th International Symposium on Computer Architecture (ISCA),
pages 63-74, Beijing, China, June 2008. [Slides \(ppt\)](#)

Another Problem due to Memory Interference

- Processors try to tolerate the latency of DRAM requests by generating multiple outstanding requests
 - Memory-Level Parallelism (MLP)
 - Out-of-order execution, non-blocking caches, runahead execution
- Effective only if the DRAM controller actually services the multiple requests in parallel in DRAM banks
- Multiple threads share the DRAM controller
- DRAM controllers are not aware of a thread's MLP
 - Can service each thread's outstanding requests **serially, not in parallel**

Bank Parallelism of a Thread



Bank access latencies of the two requests overlapped
Thread stalls for ~ONE bank access latency

Bank Parallelism Interference in DRAM

Baseline Scheduler:

2 DRAM Requests



Bank 0

Bank 1

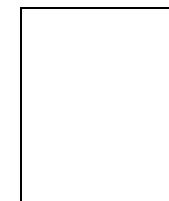
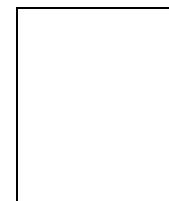
2 DRAM Requests



Bank 1

Bank 0

Bank 0 Bank 1



Thread A: Bank 0, Row 1

Thread B: Bank 1, Row 99

Thread B: Bank 0, Row 99

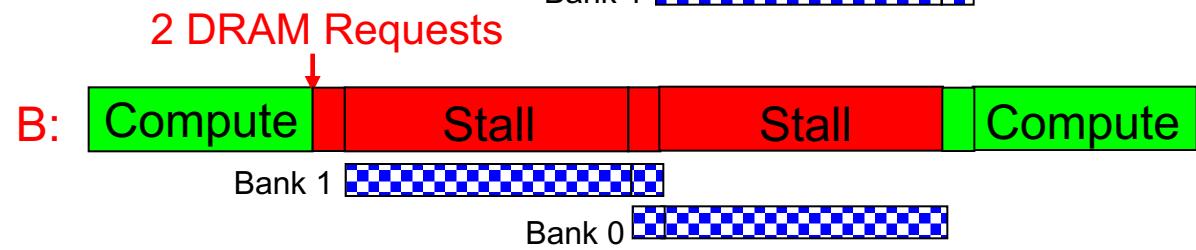
Thread A: Bank 1, Row 1

Bank access latencies of each thread serialized
Each thread stalls for ~TWO bank access latencies

Parallelism-Aware Scheduler

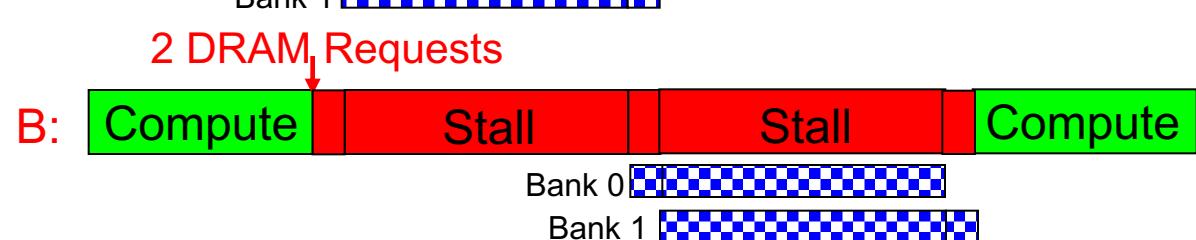
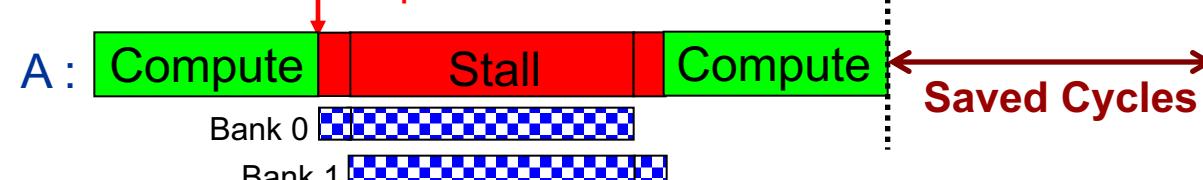
Baseline Scheduler:

2 DRAM Requests

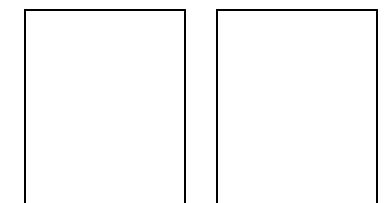


Parallelism-aware Scheduler:

2 DRAM Requests



Bank 0 Bank 1



Thread A: Bank 0, Row 1

Thread B: Bank 1, Row 99

Thread B: Bank 0, Row 99

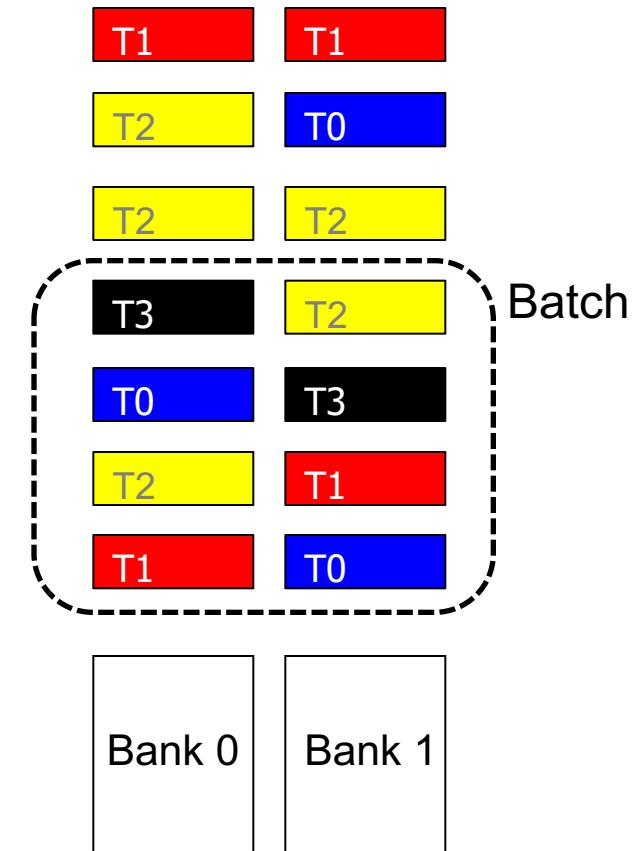
Thread A: Bank 1, Row 1

Average stall-time:
~1.5 bank access
latencies

Parallelism-Aware Batch Scheduling (PAR-BS)

■ Principle 1: Parallelism-awareness

- Schedule requests from a thread (to different banks) back to back
- Preserves each thread's bank parallelism
- But, this can cause starvation...



■ Principle 2: Request Batching

- Group a fixed number of oldest requests from each thread into a “batch”
- Service the batch before all other requests
- Form a new batch when the current one is done
- Eliminates starvation, provides fairness
- Allows parallelism-awareness within a batch

PAR-BS Components

- Request batching
- Within-batch scheduling
 - Parallelism aware

Request Batching

- Each memory request has a bit (*marked*) associated with it
- Batch formation:
 - Mark up to *Marking-Cap* oldest requests per bank for each thread
 - Marked requests constitute the batch
 - Form a new batch when no marked requests are left
- Marked requests are prioritized over unmarked ones
 - No reordering of requests across batches: **no starvation, high fairness**
- How to prioritize requests within a batch?

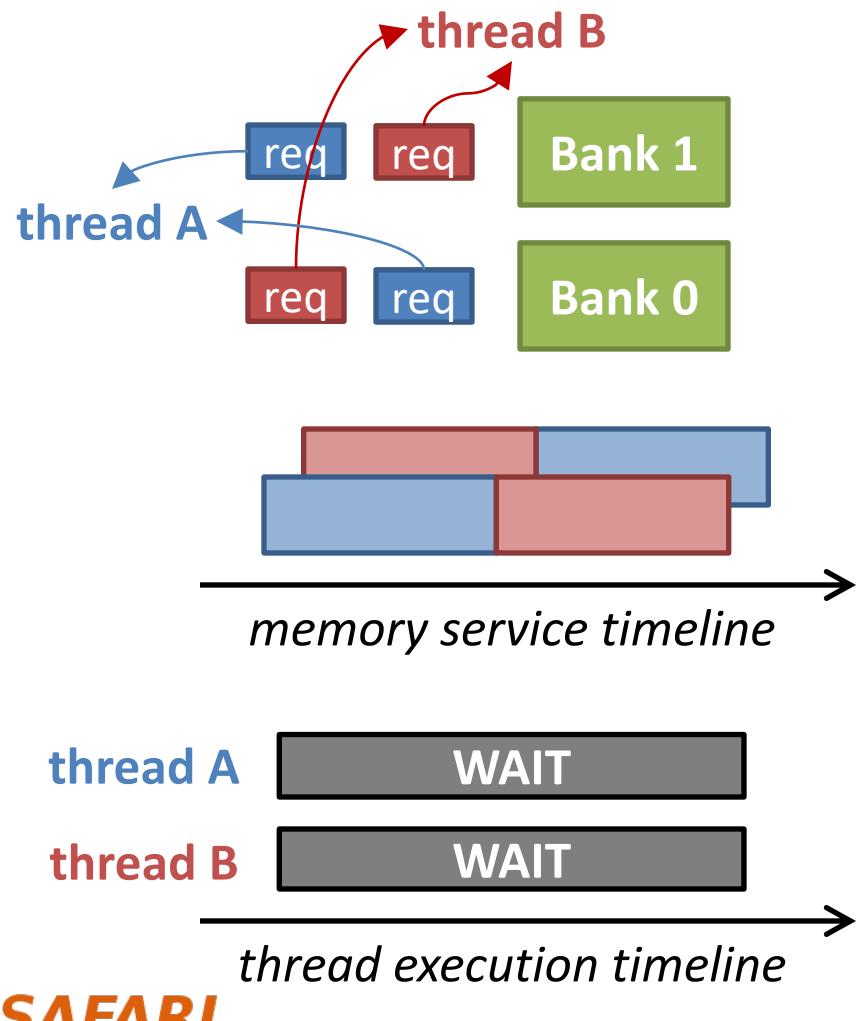
Within-Batch Scheduling

- Can use any existing DRAM scheduling policy
 - FR-FCFS (row-hit first, then oldest-first) exploits row-buffer locality
- But, we also want to preserve intra-thread bank parallelism
 - Service each thread's requests back to back

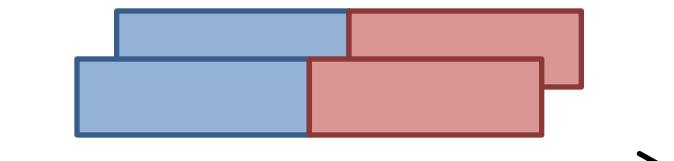
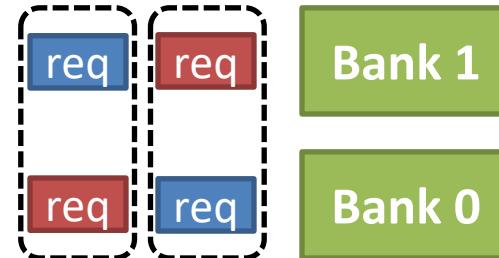
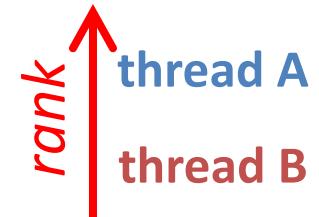
HOW?

- Scheduler **computes a ranking of threads** when the batch is formed
 - Higher-ranked threads are prioritized over lower-ranked ones
 - Improves the likelihood that requests from a thread are serviced in parallel by different banks
 - Different threads prioritized in the same order across ALL banks

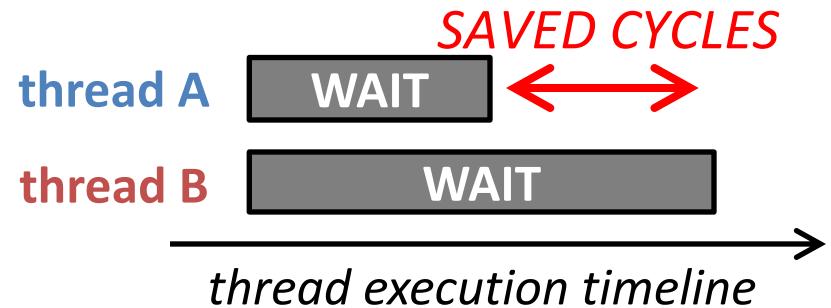
Thread Ranking



Key Idea:



memory service timeline



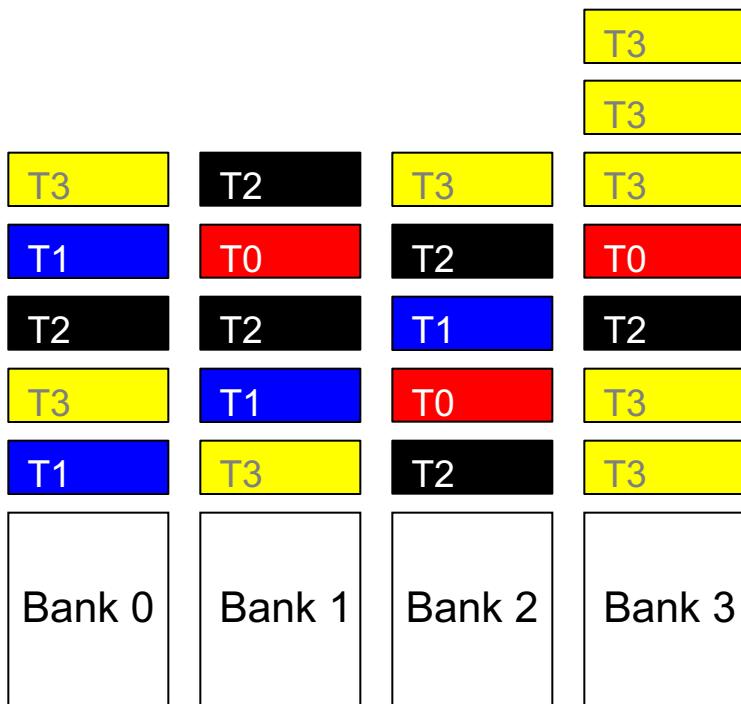
How to Rank Threads within a Batch

- Ranking scheme affects system throughput and fairness
- Maximize system throughput
 - Minimize average stall-time of threads within the batch
- Minimize unfairness (Equalize the slowdown of threads)
 - Service threads with inherently low stall-time early in the batch
 - Insight: delaying memory non-intensive threads results in high slowdown
- Shortest stall-time first (shortest job first) ranking
 - Provides optimal system throughput [Smith, 1956]*
 - Controller estimates each thread's stall-time within the batch
 - Ranks threads with shorter stall-time higher

* W.E. Smith, “Various optimizers for single stage production,” Naval Research Logistics Quarterly, 1956.

Shortest Stall-Time First Ranking

- Maximum number of marked requests to any bank (max-bank-load)
 - Rank thread with lower max-bank-load higher (~ low stall-time)
- Total number of marked requests (total-load)
 - Breaks ties: rank thread with lower total-load higher

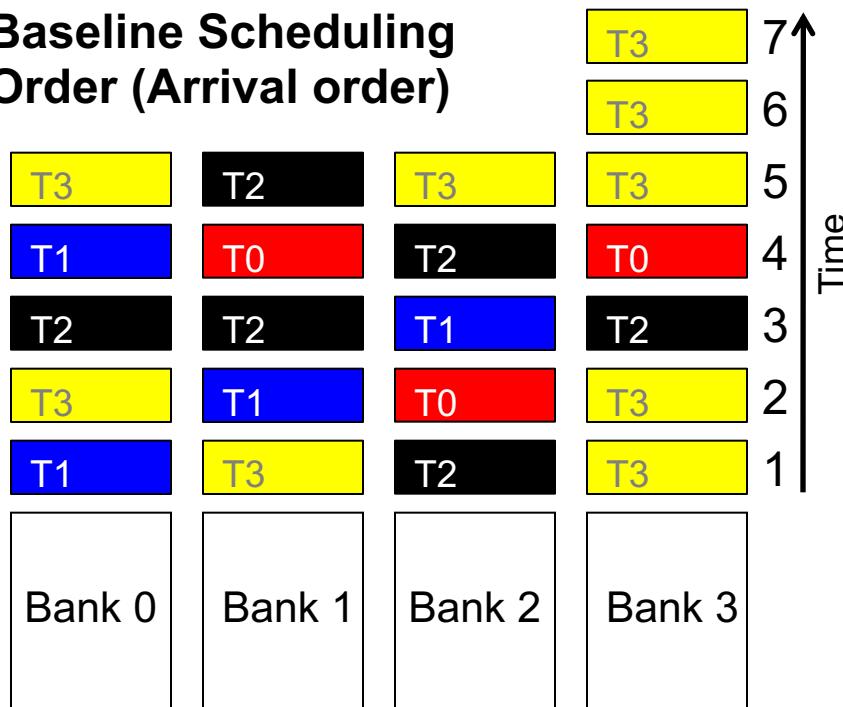


	max-bank-load	total-load

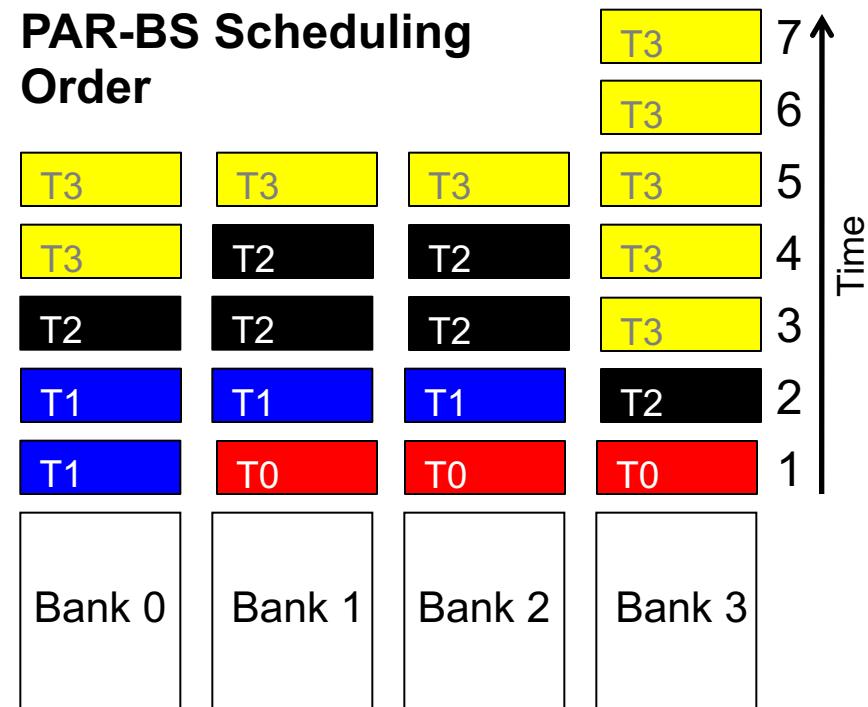
Ranking:
T0 > T1 > T2 > T3

Example Within-Batch Scheduling Order

Baseline Scheduling Order (Arrival order)



PAR-BS Scheduling Order



Ranking: $T_0 > T_1 > T_2 > T_3$

	T0	T1	T2	T3
Stall times				

AVG: 5 bank access latencies

	T0	T1	T2	T3
Stall times				

AVG: 3.5 bank access latencies

Putting It Together: PAR-BS Scheduling Policy

■ PAR-BS Scheduling Policy

(1) Marked requests first

Batching

(2) Row-hit requests first

Parallelism-aware
within-batch
scheduling

(3) Higher-rank thread first (shortest stall-time first)

(4) Oldest first

■ Three properties:

- Exploits row-buffer locality **and** intra-thread bank parallelism
- Work-conserving
 - Services unmarked requests to banks without marked requests
- Marking-Cap is important
 - Too small cap: destroys row-buffer locality
 - Too large cap: penalizes memory non-intensive threads

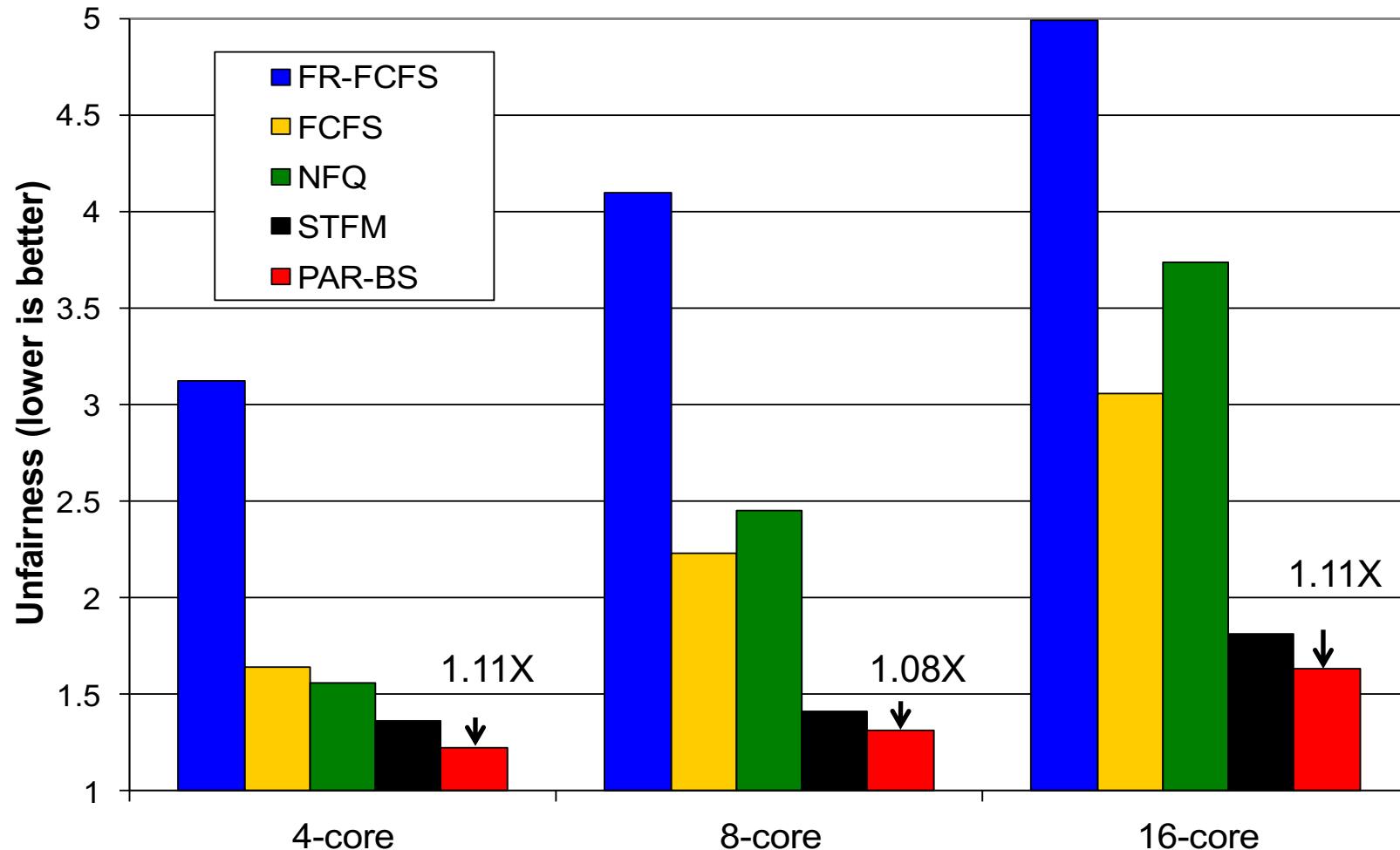
■ Many more trade-offs analyzed in the paper

Hardware Cost

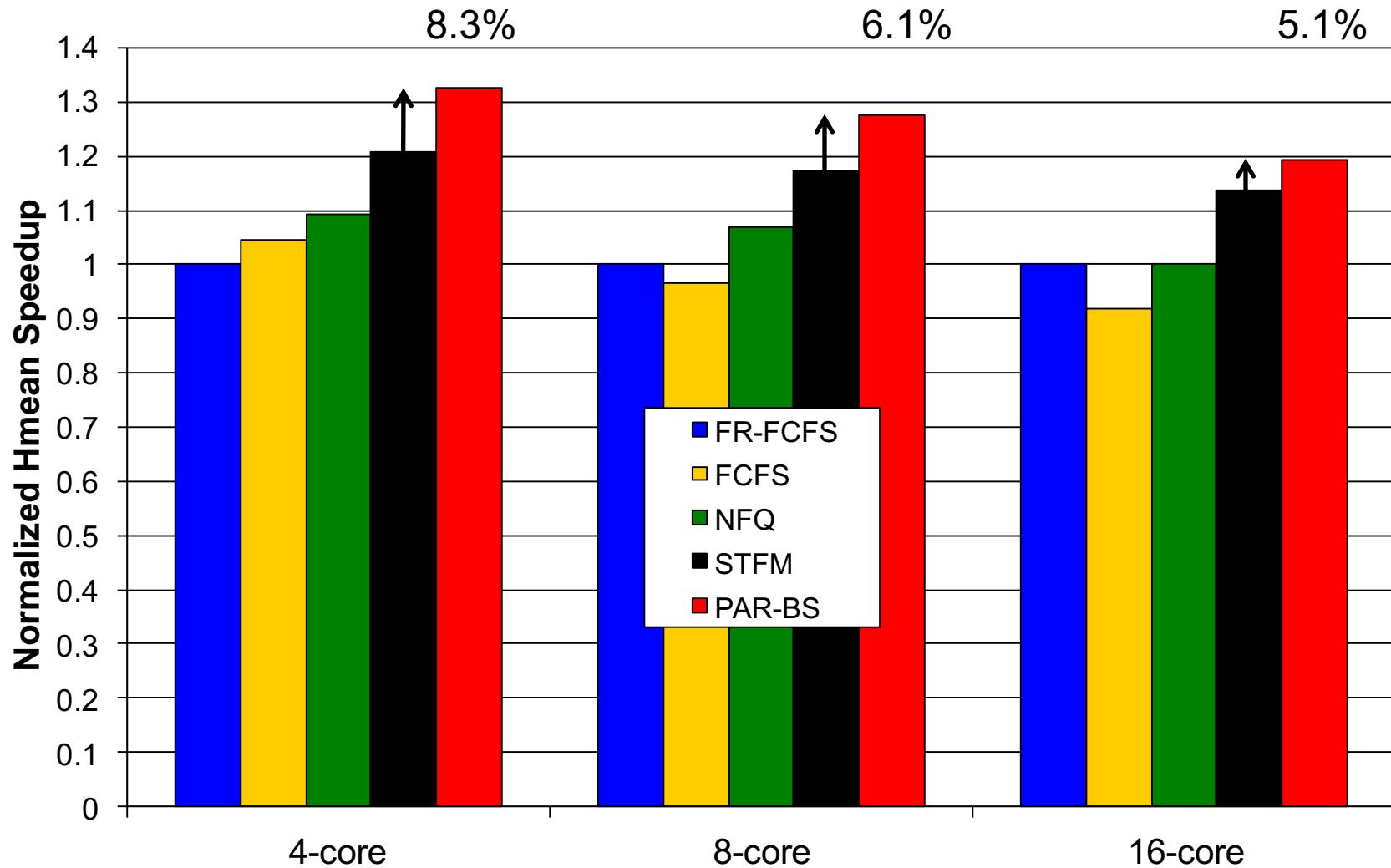
- <1.5KB storage cost for
 - 8-core system with 128-entry memory request buffer
- No complex operations (e.g., divisions)
- Not on the critical path
 - Scheduler makes a decision only every DRAM cycle

Unfairness on 4-, 8-, 16-core Systems

Unfairness = MAX Memory Slowdown / MIN Memory Slowdown [MICRO 2007]



System Performance (Hmean-speedup)



PAR-BS Pros and Cons

- Upsides:
 - First scheduler to address bank parallelism destruction across multiple threads
 - Simple mechanism (vs. STFM)
 - Batching provides fairness
 - Ranking enables parallelism awareness

- Downsides:
 - Does not always prioritize the latency-sensitive applications

More on PAR-BS

- 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)]
One of the 12 computer architecture papers of 2008 selected as Top Picks by IEEE Micro.

Parallelism-Aware Batch Scheduling:

Enhancing both Performance and Fairness of Shared DRAM Systems

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

More on PAR-BS

- Onur Mutlu and Thomas Moscibroda,

"Parallelism-Aware Batch Scheduling: Enabling High-Performance and Fair Memory Controllers"

*IEEE Micro, Special Issue: Micro's Top Picks from 2008 Computer Architecture Conferences (**MICRO TOP PICKS**)*, Vol. 29, No. 1, pages 22-32, January/February 2009.

PARALLELISM-AWARE BATCH SCHEDULING: ENABLING HIGH-PERFORMANCE AND FAIR SHARED MEMORY CONTROLLERS

UNCONTROLLED INTERTHREAD INTERFERENCE IN MAIN MEMORY CAN DESTROY INDIVIDUAL THREADS' MEMORY-LEVEL PARALLELISM, EFFECTIVELY SERIALIZING THE MEMORY REQUESTS OF A THREAD WHOSE LATENCIES WOULD OTHERWISE HAVE LARGELY OVERLAPPED, THEREBY REDUCING SINGLE-THREAD PERFORMANCE. THE PARALLELISM-AWARE BATCH SCHEDULER PRESERVES EACH THREAD'S MEMORY-LEVEL PARALLELISM, ENSURES FAIRNESS AND STARVATION FREEDOM, AND SUPPORTS SYSTEM-LEVEL THREAD PRIORITIES.

ATLAS Memory Scheduler

Yoongu Kim, Dongsu Han, Onur Mutlu, and Mor Harchol-Balter,

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

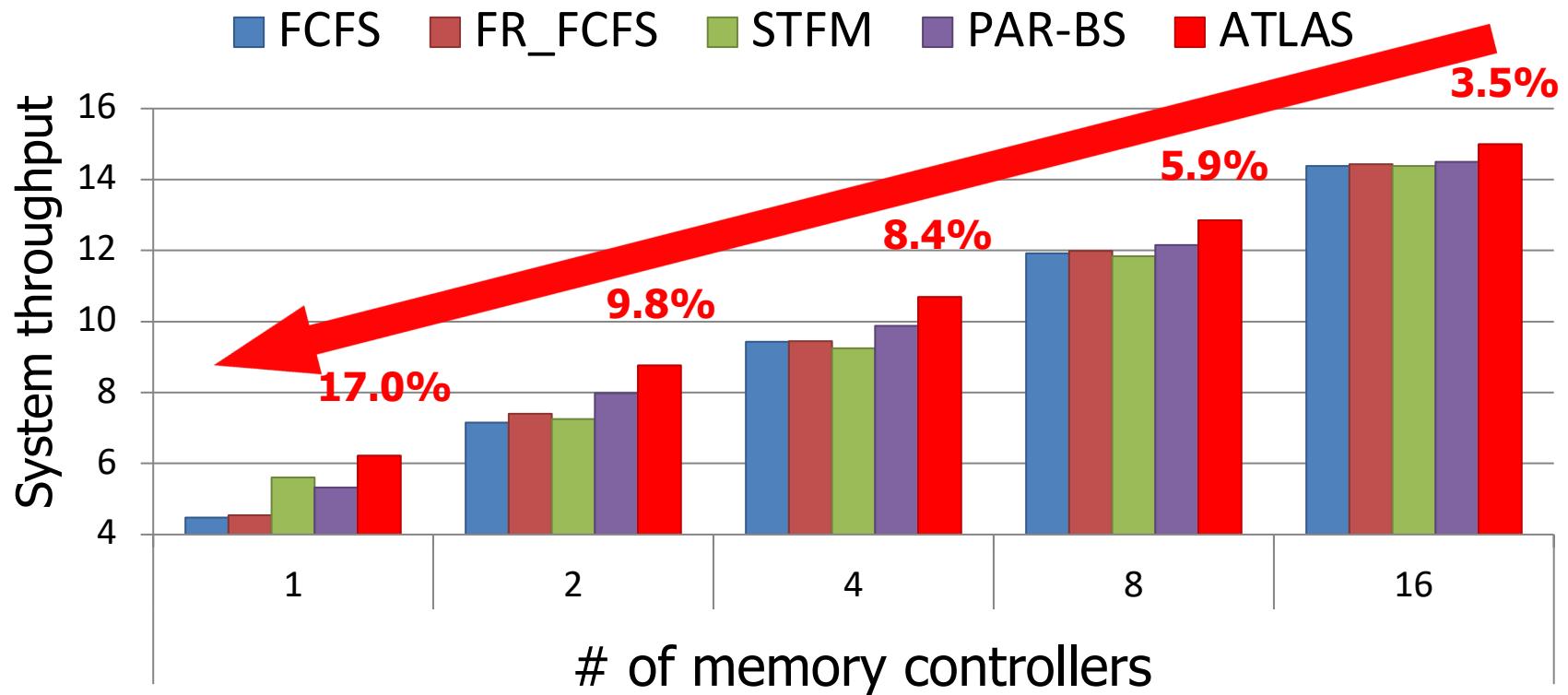
*16th International Symposium on High-Performance Computer Architecture (HPCA),
Bangalore, India, January 2010.* [Slides \(pptx\)](#)

ATLAS: Summary

- Goal: To maximize system performance
- Main idea: Prioritize the thread that has attained the least service from the memory controllers (Adaptive per-Thread Least Attained Service Scheduling)
 - Rank threads based on attained service in the past time interval(s)
 - Enforce thread ranking in the memory scheduler during the current interval
- Why it works: Prioritizes “light” (memory non-intensive) threads that are more likely to keep their cores busy

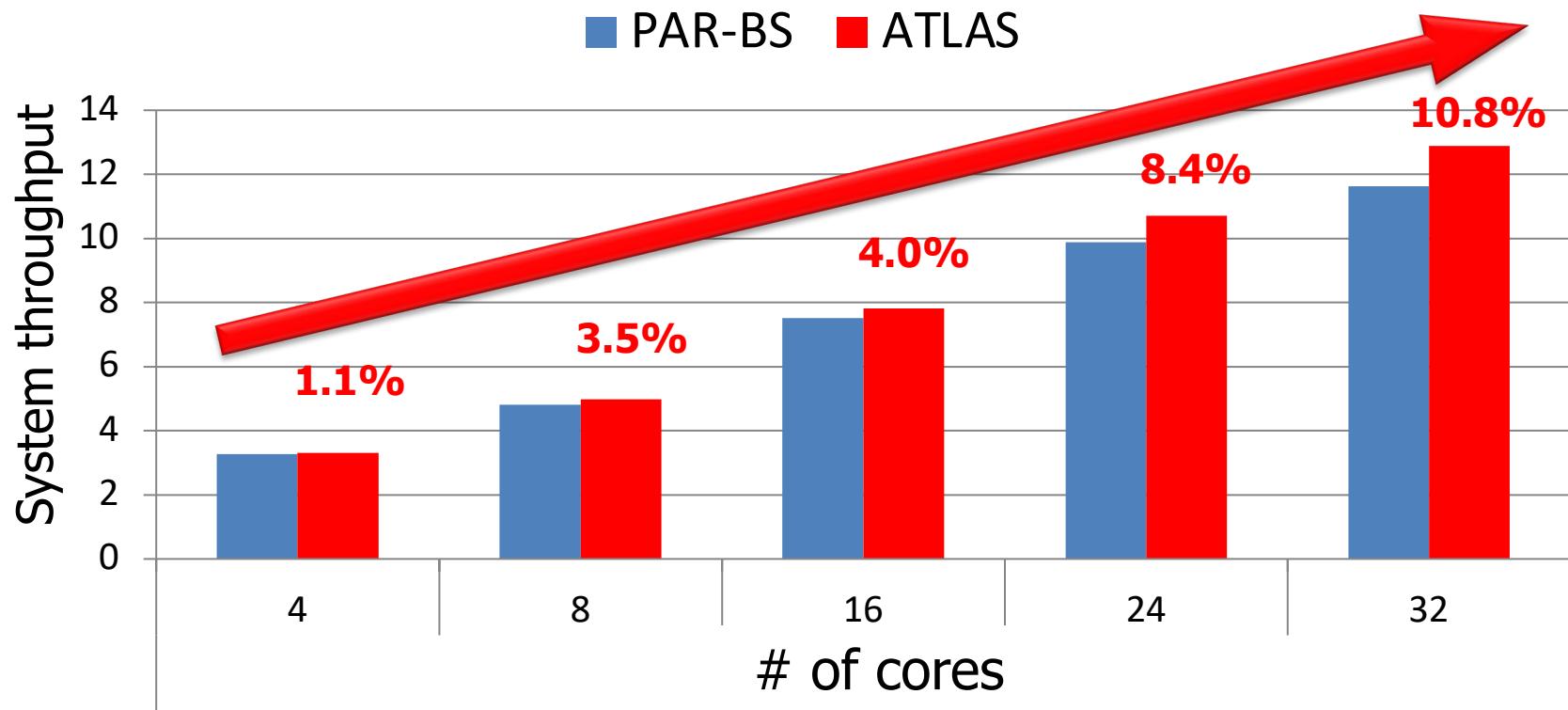
System Throughput: 24-Core System

System throughput = \sum Speedup



ATLAS consistently provides higher system throughput than all previous scheduling algorithms

System Throughput: 4-MC System



of cores increases → ATLAS performance benefit increases

ATLAS Pros and Cons

- Upsides:
 - Good at improving overall throughput (compute-intensive threads are prioritized)
 - Low complexity
 - Coordination among controllers happens infrequently

- Downsides:
 - Lowest/medium ranked threads get delayed significantly → high unfairness

More on ATLAS Memory Scheduler

- 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)

Best paper session. One of the four papers nominated for the Best Paper Award by the Program Committee.

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

Yoongu Kim Dongsu Han Onur Mutlu Mor Harchol-Balter

Carnegie Mellon University

Computer Architecture

Lecture 12: Memory Controllers

Ataberk Olgun

Prof. Onur Mutlu

ETH Zürich

Fall 2023

3 November 2023