# Memory Systems
## and Memory-Centric Computing Systems

## Lecture 2a: Memory Controllers

Prof. Onur Mutlu

omutlu@gmail.com

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

13 June 2019

TU Wien Fast Course 2019

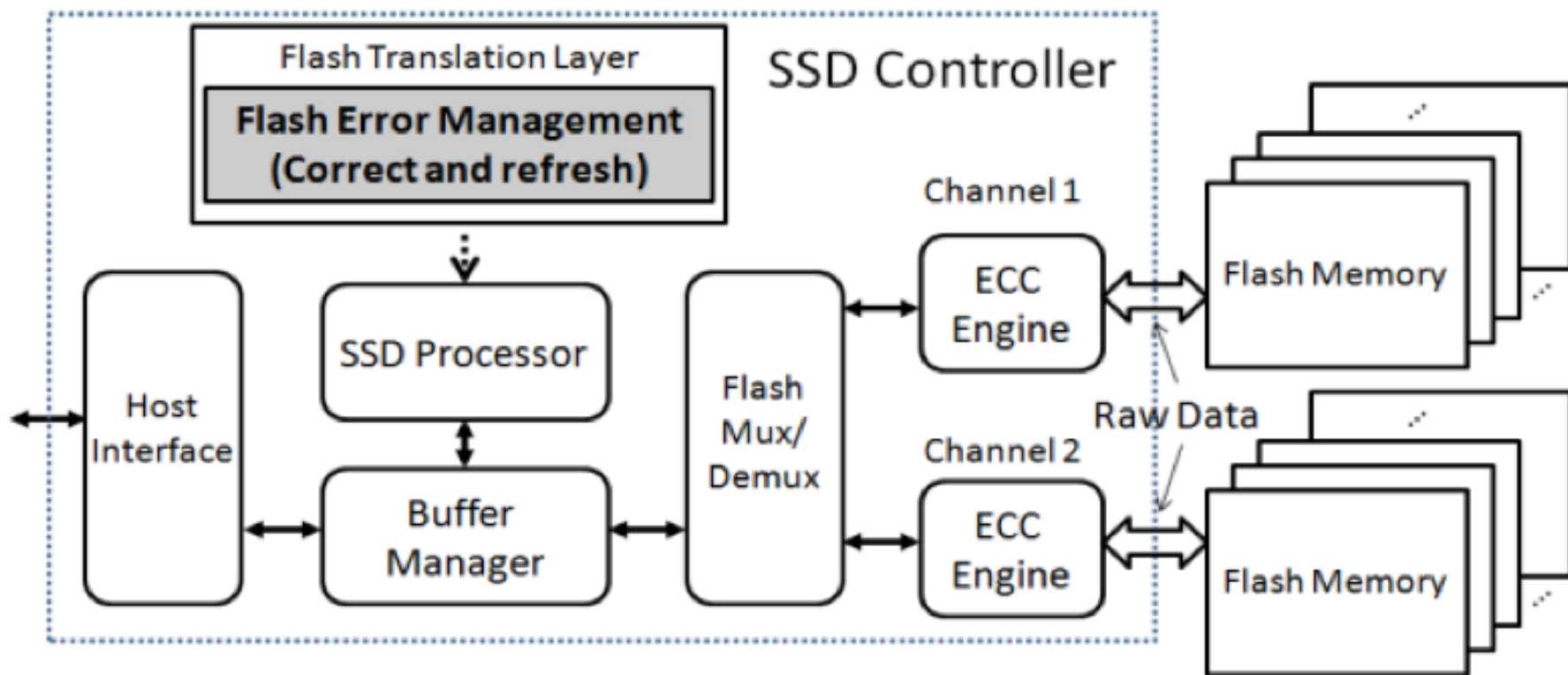**SAFARI**    **ETH** *zürich*    **Carnegie Mellon**

# Memory Controllers

# DRAM versus Other Types of Memories

- Long latency memories have similar characteristics that need to be controlled.

- The following discussion 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 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: error correction, garbage collection, page remapping, …



Cai+, "Flash Correct-and-Refresh: Retention-Aware Error Management for Increased Flash Memory Lifetime", ICCD 2012.

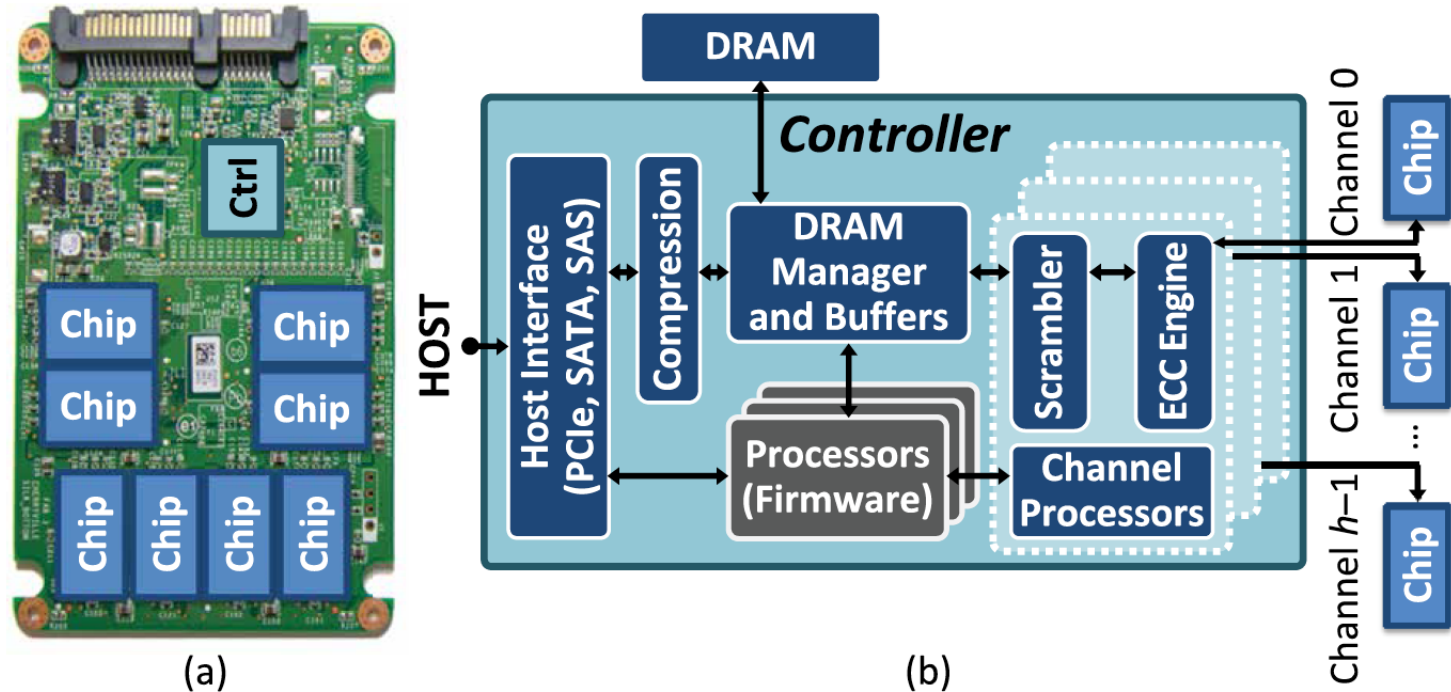# Another View of the SSD Controller



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

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

# On Modern SSD Controllers (I)

INVITED PAPER

# 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

https://arxiv.org/pdf/1706.08642

6

# 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*

# DRAM Types

- DRAM has different types with different interfaces optimized for different purposes
  - Commodity: DDR, DDR2, DDR3, DDR4, …
  - Low power (for mobile): LPDDR1, …, LPDDR5, …
  - High bandwidth (for graphics): GDDR2, …, GDDR5, …
  - Low latency: eDRAM, RLDRAM, …
  - 3D stacked: WIO, HBM, HMC, …
  - …
- 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)

# DRAM Types (circa 2015)

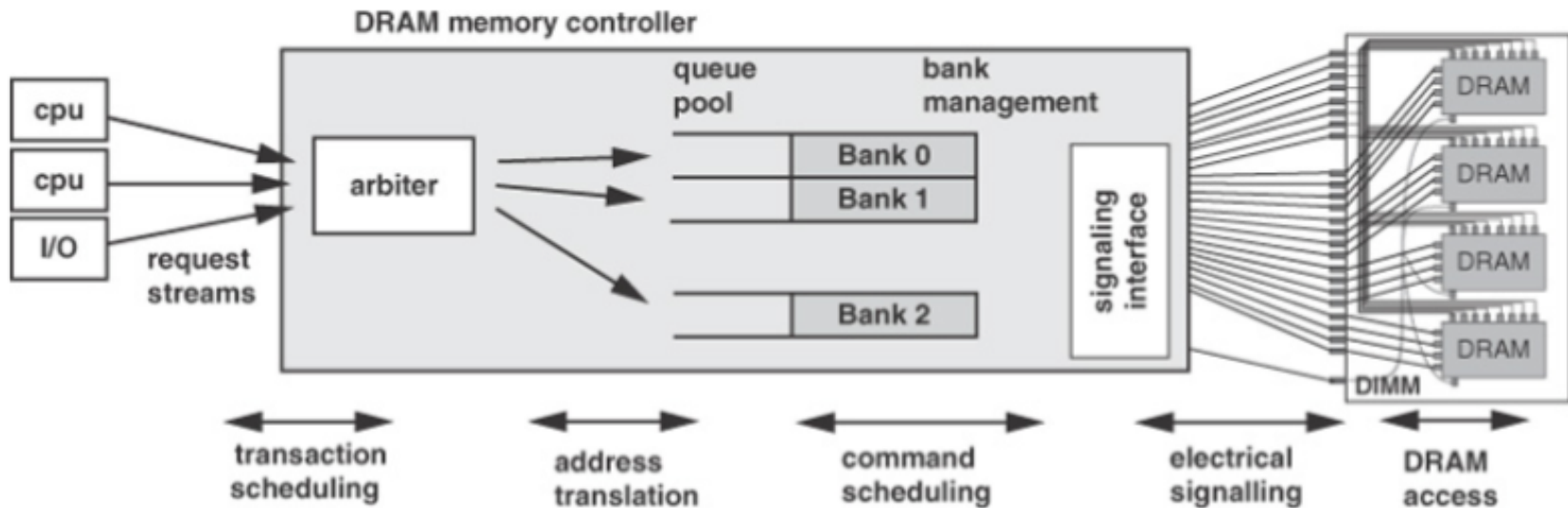| Segment | DRAM Standards & Architectures |
|---|---|
| Commodity | DDR3 (2007) [14]; DDR4 (2012) [18] |
| Low-Power | LPDDR3 (2012) [17]; LPDDR4 (2014) [20] |
| Graphics | GDDR5 (2009) [15] |
| Performance | eDRAM [28], [32]; RLDRAM3 (2011) [29] |
| 3D-Stacked | WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (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 et al., "Ramulator: A Fast and Extensible DRAM Simulator," IEEE Comp Arch Letters 2015.
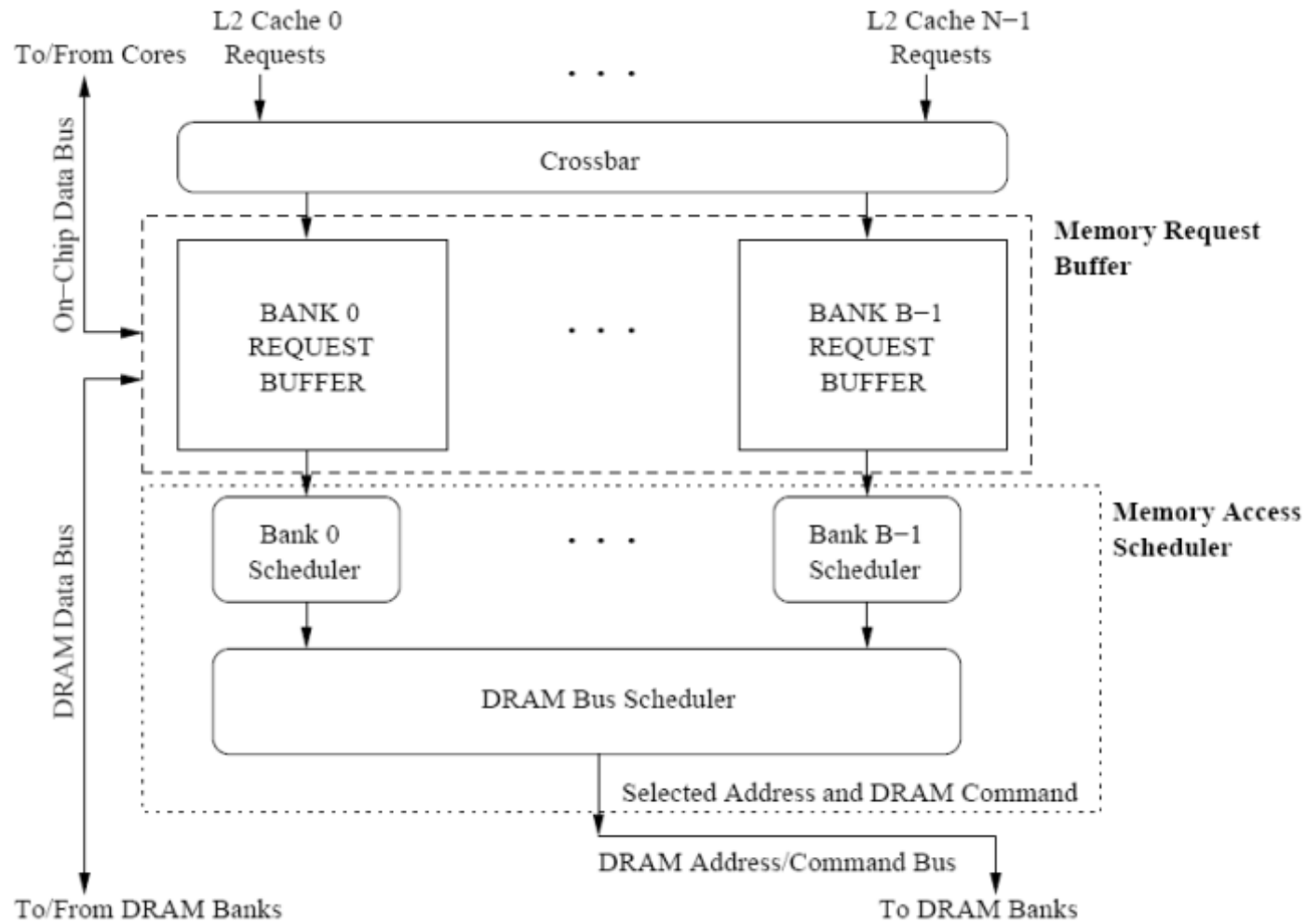
# DRAM Controller: Functions

- 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 to 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

Mutlu+, "Stall-Time Fair Memory Scheduling," MICRO 2007.

# DRAM Scheduling Policies (I)

- **FCFS** (first come first served)
  - Oldest request first

- **FR-FCFS** (first ready, first come first served)
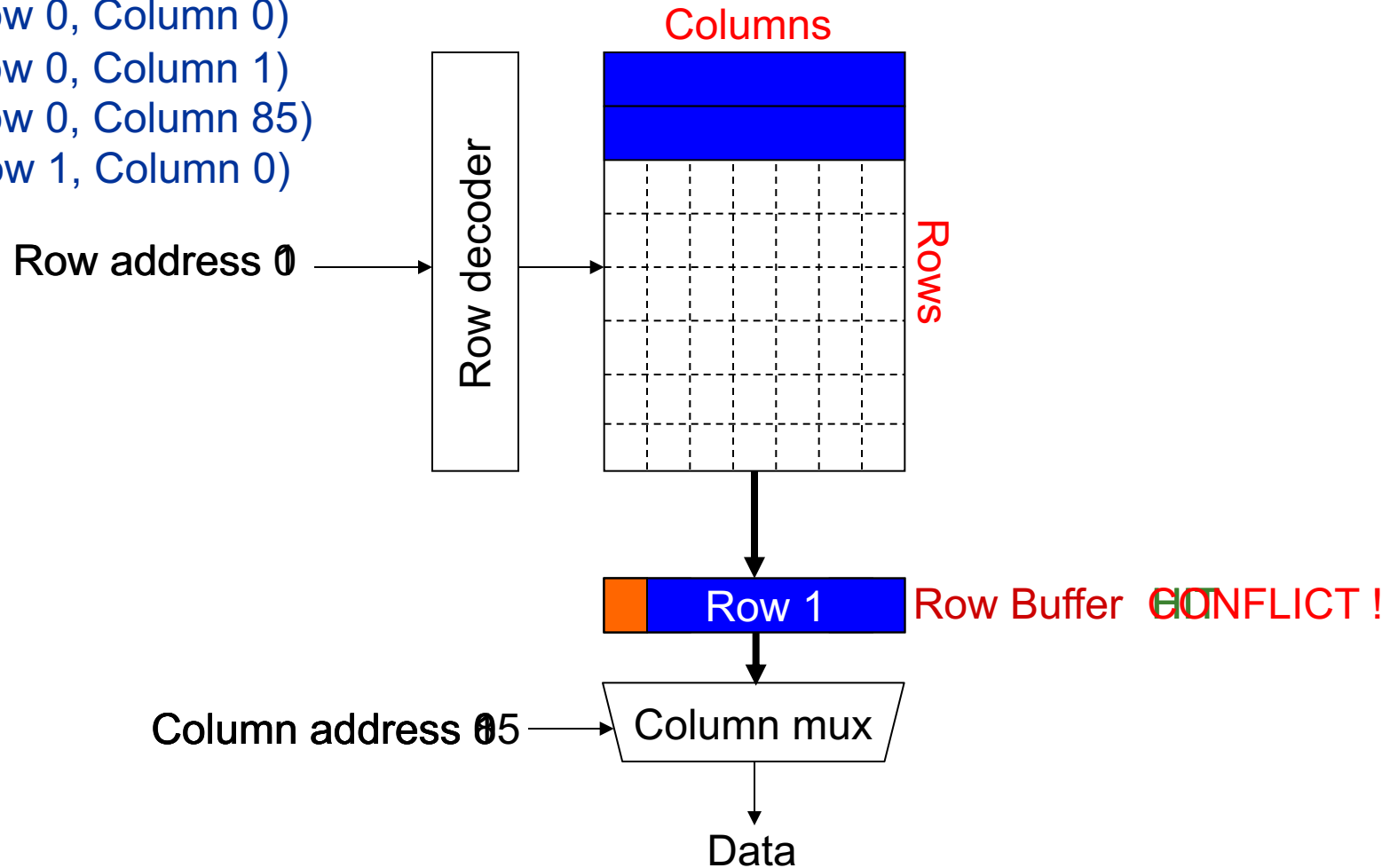  1. Row-hit first
  2. Oldest first

  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, older commands prioritized over younger ones

# Review: DRAM Bank Operation

Access Address:
(Row 0, Column 0)
(Row 0, Column 1)
(Row 0, Column 85)
(Row 1, Column 0)

Columns

Row decoder

Rows

Row address 01

Row 1    Row Buffer   CONFLICT !

Column address 085    Column mux

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

- 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 | $^tRP$ | 11 | Activate to read/write | $^tRCD$ | 11 |
| Read column address strobe | $CL$ | 11 | Write column address strobe | $CWL$ | 8 |
| Additive | $AL$ | 0 | Activate to activate | $^tRC$ | 39 |
| Activate to precharge | $^tRAS$ | 28 | Read to precharge | $^tRTP$ | 6 |
| Burst length | $^tBL$ | 4 | Column address strobe to column address strobe | $^tCCD$ | 4 |
| Activate to activate (different bank) | $^tRRD$ | 6 | Four activate windows | $^tFAW$ | 24 |
| Write to read | $^tWTR$ | 6 | Write recovery | $^tWR$ | 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.
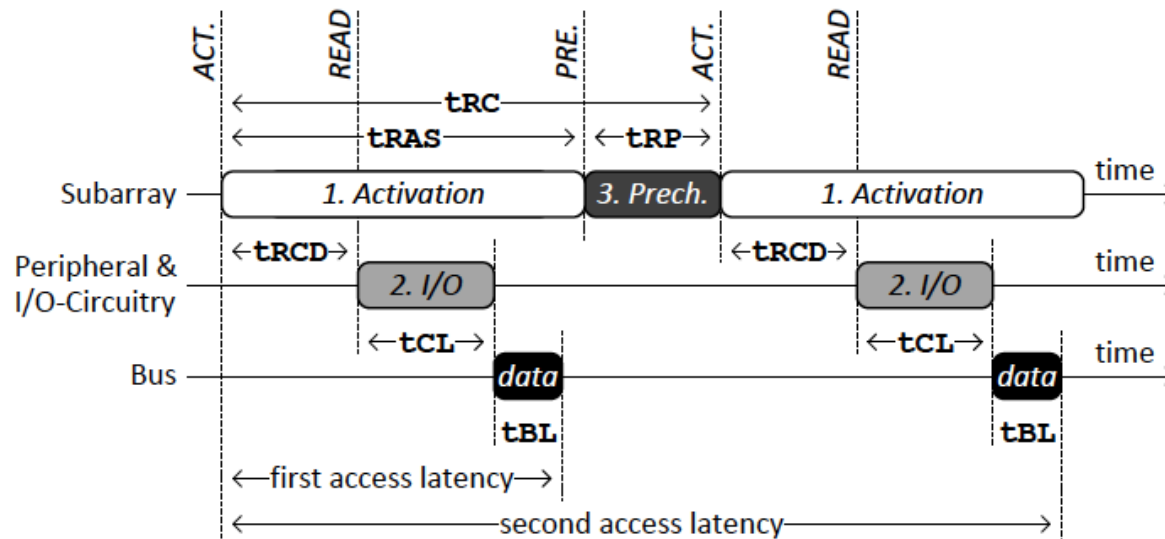


Figure 5. Three Phases of DRAM Access

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

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

# 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↔RowCmd | | | RowCmd↔ColCmd | | | ColCmd↔ColCmd | | | ColCmd→DATA | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Name | $tRC$ | $tRAS$ | $tRP$ | $tRCD$ | $tRTP$ | $tWR^*$ | $tCCD$ | $tRTW^\dagger$ | $tWTR^*$ | $CL$ | $CWL$ |
| Commands | A→A | A→P | P→A | A→R/W | R→P | W*→P | R(W)→R(W) | R→W | W*→R | R→DATA | W→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.
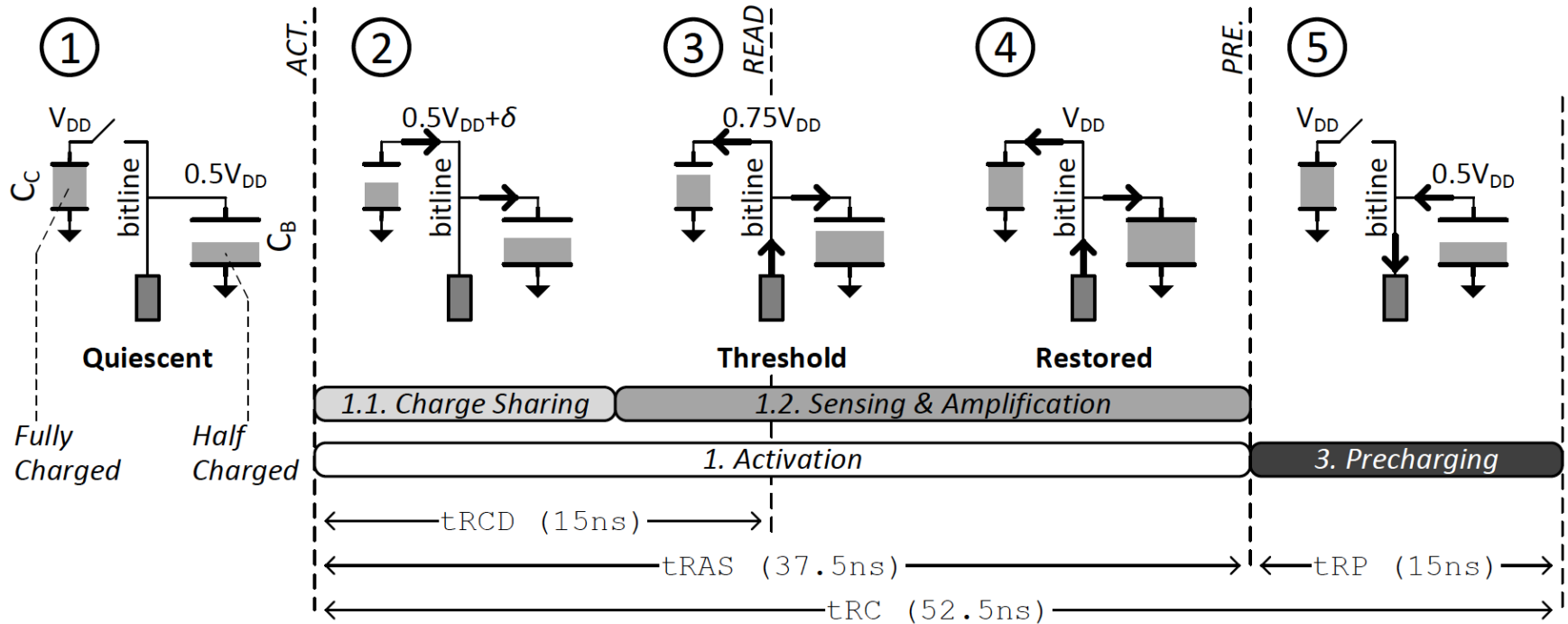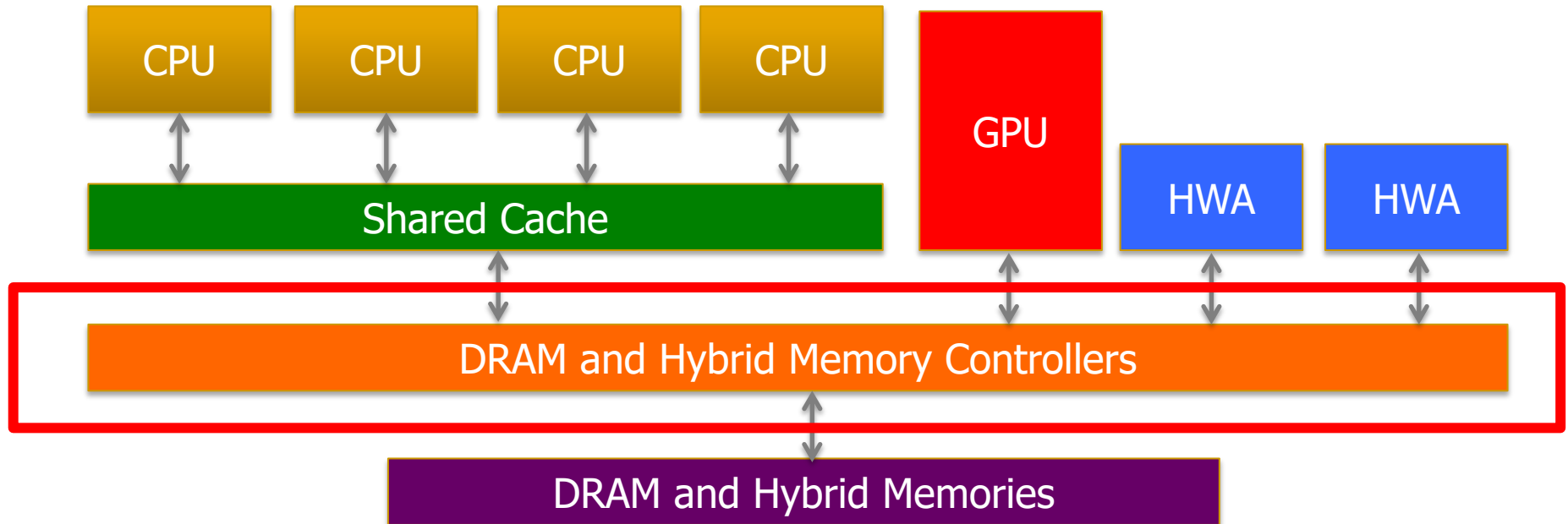
# Why So Many Timing Constraints? (II)



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<br>ACT → WRITE | tRCD | 15ns |
| | ACT → PRE | tRAS | 37.5ns |
| 2 | READ → data<br>WRITE → data | tCL<br>tCWL | 15ns<br>11.25ns |
| | data burst | tBL | 7.5ns |
| 3 | PRE → ACT | tRP | 15ns |
| 1 & 3 | ACT → ACT | tRC<br>(tRAS+tRP) | 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, …

# 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?

# 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.

Ipek+, "Self Optimizing Memory Controllers: A Reinforcement Learning Approach," ISCA 2008.

# Self-Optimizing DRAM Controllers



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

**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

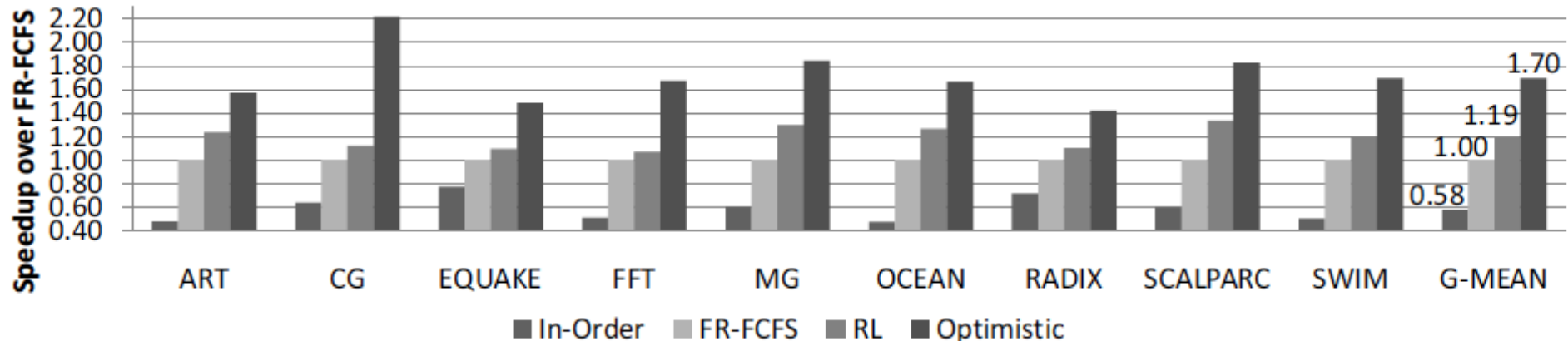- Read - store miss

- 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
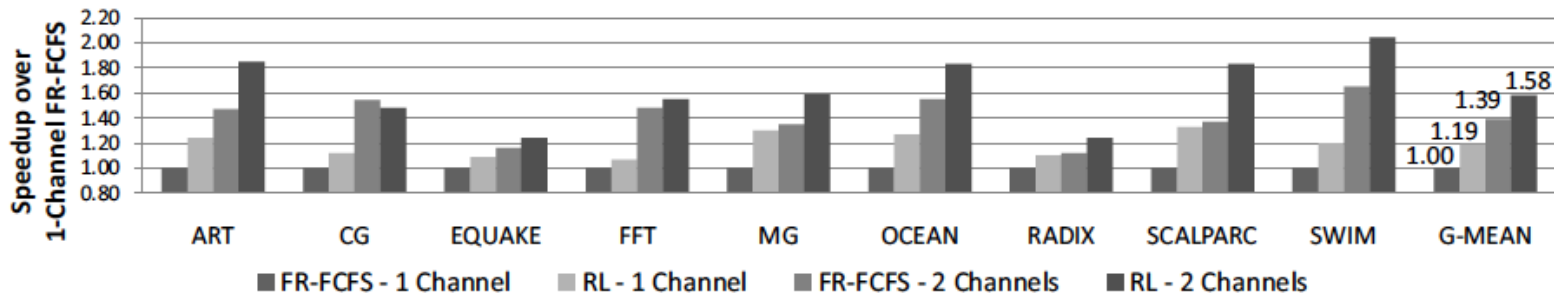


Figure 15: Performance comparison of FR-FCFS and RL-based memory controllers on systems with 6.4GB/s and 12.8GB/s peak DRAM bandwidth

# 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 Mutlu[2]    José F. Martínez[1]    Rich Caruana[1]

[1] Cornell 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)

# Memory Interference

# 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

# Uncontrolled Interference: An Example



Multi-Core Chip

stream

random

unfairness

L2 CACHE

L2 CACHE

INTERCONNECT

DRAM MEMORY CONTROLLER

Shared DRAM Memory System

DRAM Bank 0

DRAM Bank 1

DRAM Bank 2

DRAM Bank 3

# A Memory Performance Hog

```
// initialize large arrays A, B

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

```
// initialize large arrays A, B

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

**STREAM**

**RANDOM**

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

- 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?



Memory Request Buffer

T0: Row 0
T1: Row 5 / T0: Row 0
T1: Row 111 / T0: Row 0
T1: Row 16 / T0: Row 0

Row decoder

Row Buffer

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.
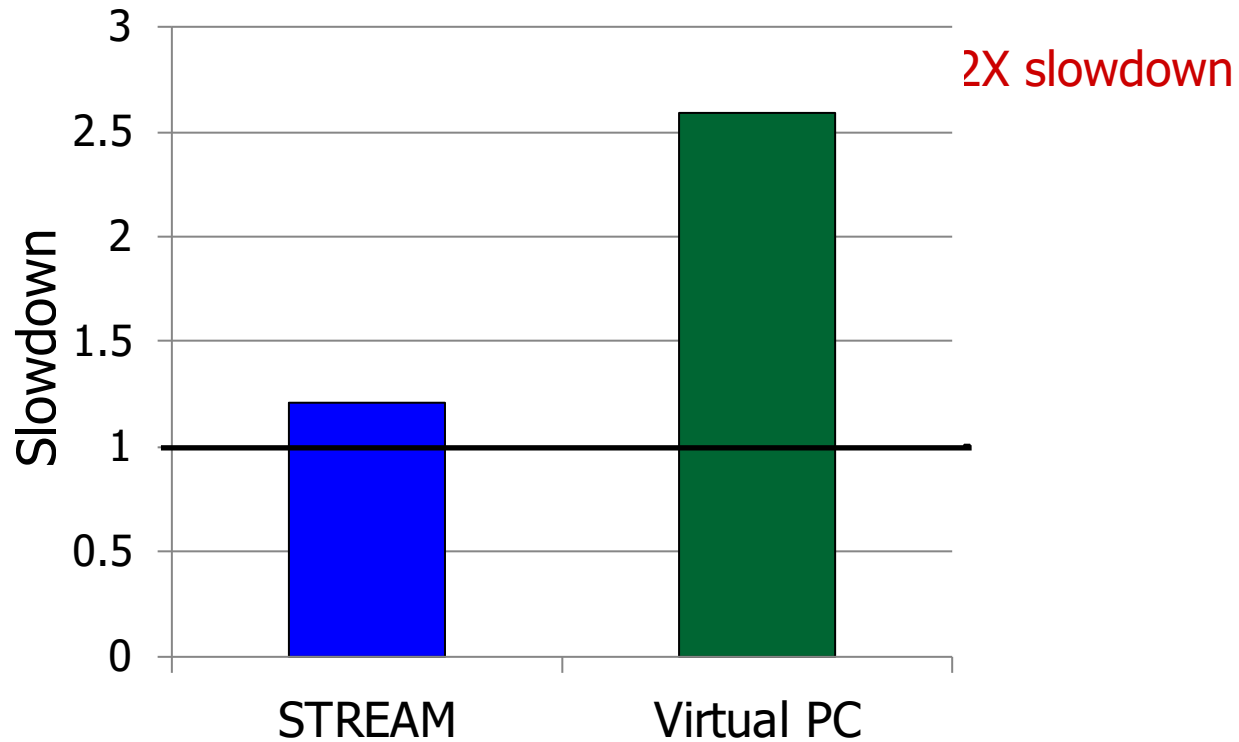
# Unfair Slowdowns due to Interference

Moscibroda and Mutlu, "Memory performance attacks: Denial of memory service in multi-core systems," 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)

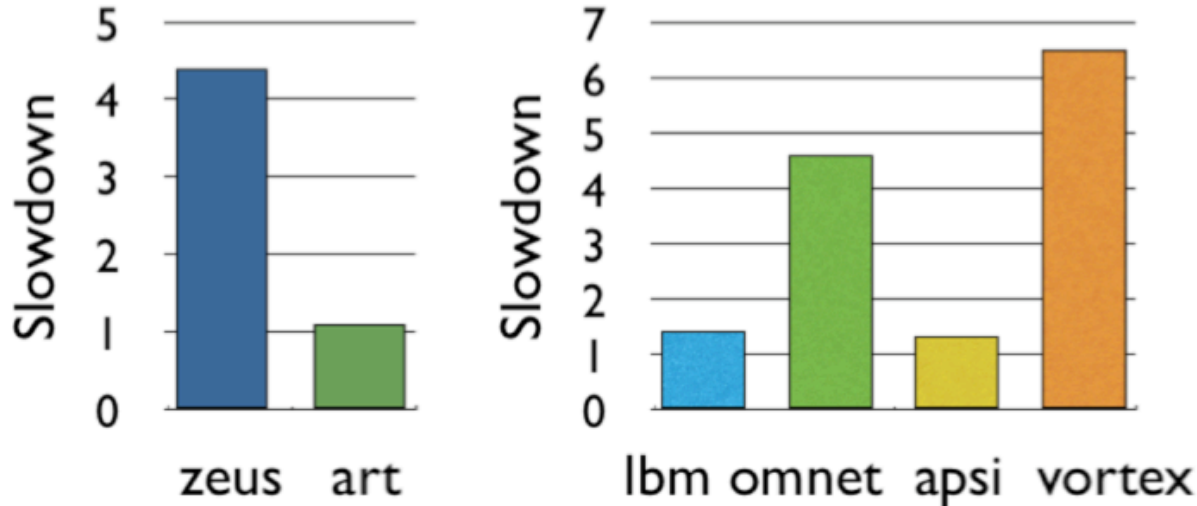Moscibroda and Mutlu, "Memory Performance Attacks," USENIX Security 2007.

# 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**

# More on Memory Performance Attacks

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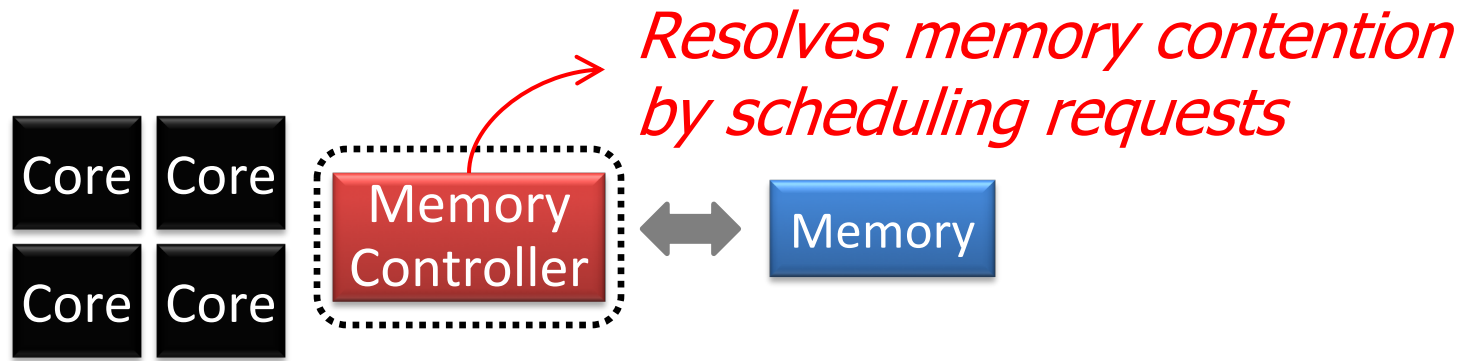  *Proceedings of the 16th USENIX Security Symposium* (**USENIX SECURITY**), pages 257-274, Boston, MA, August 2007. Slides (ppt)

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

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

# 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 Scheduling

Core Core
Core Core

Memory Controller ⬌ Memory

*Resolves memory contention by scheduling requests*

- 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

**SAFARI**

# QoS-Aware Memory: Readings (I)

- 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

# QoS-Aware Memory: Readings (II)

- 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

# QoS-Aware Memory: Readings (III)

- 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

# QoS-Aware Memory: Readings (IV)

- 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

# QoS-Aware Memory: Readings (V)

- 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

# QoS-Aware Memory: Readings (VI)

- Lavanya Subramanian, Donghyuk Lee, Vivek Seshadri, Harsha Rastogi, and Onur Mutlu,
  **"BLISS: Balancing Performance, Fairness and Complexity in Memory Access Scheduling"**
  *IEEE Transactions on Parallel and Distributed Systems* (**TPDS**), to appear in 2016.  arXiv.org version, April 2015.
  An earlier version as *SAFARI Technical Report*, TR-SAFARI-2015-004, Carnegie Mellon University, March 2015.
  [Source Code]

# BLISS: Balancing Performance, Fairness and Complexity in Memory Access Scheduling

Lavanya Subramanian, Donghyuk Lee, Vivek Seshadri, Harsha Rastogi, and Onur Mutlu

# QoS-Aware Memory: Readings (VII)

- 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[†]   Kevin Kai-Wei Chang[†]   Lavanya Subramanian[†]   Gabriel H. Loh[‡]   Onur Mutlu[†]

[†]Carnegie Mellon University
{rachata,kevincha,lsubrama,onur}@cmu.edu

[‡]Advanced Micro Devices, Inc.
gabe.loh@amd.com

# QoS-Aware Memory: Readings (VIII)

- 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

# QoS-Aware Memory: Readings (IX)

- 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

# QoS-Aware Memory: Readings (X)

- 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

# Some More Suggested Readings

# Some Key Readings on DRAM (I)

- **DRAM Organization and Operation**

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

    https://people.inf.ethz.ch/omutlu/pub/tldram_hpca13.pdf

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

    https://people.inf.ethz.ch/omutlu/pub/salp-dram_isca12.pdf

  - Lee et al., "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost," ACM TACO 2016.

    https://people.inf.ethz.ch/omutlu/pub/smla_high-bandwidth-3d-stacked-memory_taco16.pdf

# Some Key Readings on DRAM (II)

- **DRAM Refresh**

  - Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
    https://people.inf.ethz.ch/omutlu/pub/raidr-dram-refresh_isca12.pdf

  - Chang et al., "Improving DRAM Performance by Parallelizing Refreshes with Accesses," HPCA 2014.
    https://people.inf.ethz.ch/omutlu/pub/dram-access-refresh-parallelization_hpca14.pdf

  - Patel et al., "The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions," ISCA 2017.
    https://people.inf.ethz.ch/omutlu/pub/reaper-dram-retention-profiling-lpddr4_isca17.pdf

# Reading on Simulating Main Memory

- How to evaluate future main memory systems?
- An open-source simulator and its brief description

- Yoongu Kim, Weikun Yang, and Onur Mutlu,
  **"Ramulator: A Fast and Extensible DRAM Simulator"**
  *IEEE Computer Architecture Letters* (**CAL**), March 2015.
  [Source Code]

# Some Key Readings on Memory Control 1

❑ Mutlu+, "Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems," ISCA 2008.

https://people.inf.ethz.ch/omutlu/pub/parbs_isca08.pdf

❑ Kim et al., "Thread Cluster Memory Scheduling: Exploiting Differences in Memory Access Behavior," MICRO 2010.

https://people.inf.ethz.ch/omutlu/pub/tcm_micro10.pdf

❑ Subramanian et al., "BLISS: Balancing Performance, Fairness and Complexity in Memory Access Scheduling," TPDS 2016.

https://people.inf.ethz.ch/omutlu/pub/bliss-memory-scheduler_ieee-tpds16.pdf

❑ Usui et al., "DASH: Deadline-Aware High-Performance Memory Scheduler for Heterogeneous Systems with Hardware Accelerators," TACO 2016.

https://people.inf.ethz.ch/omutlu/pub/dash_deadline-aware-heterogeneous-memory-scheduler_taco16.pdf

# Some Key Readings on Memory Control 2

- Ipek+, "Self Optimizing Memory Controllers: A Reinforcement Learning Approach," ISCA 2008.

  https://people.inf.ethz.ch/omutlu/pub/rlmc_isca08.pdf

- Ebrahimi et al., "Fairness via Source Throttling: A Configurable and High-Performance Fairness Substrate for Multi-Core Memory Systems," ASPLOS 2010.

  https://people.inf.ethz.ch/omutlu/pub/fst_asplos10.pdf

- Subramanian et al., "The Application Slowdown Model: Quantifying and Controlling the Impact of Inter-Application Interference at Shared Caches and Main Memory," MICRO 2015.

  https://people.inf.ethz.ch/omutlu/pub/application-slowdown-model_micro15.pdf

- Lee et al., "Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM," PACT 2015.

  https://people.inf.ethz.ch/omutlu/pub/decoupled-dma_pact15.pdf

SAFARI

# More Readings

- To come as we cover the future topics

- Search for "DRAM" or "Memory" in:
  - https://people.inf.ethz.ch/omutlu/projects.htm

# Optional Slides:
## Inside A DRAM Chip

# DRAM Module and Chip

# Goals

- Cost

- Latency

- Bandwidth

- Parallelism

- Power
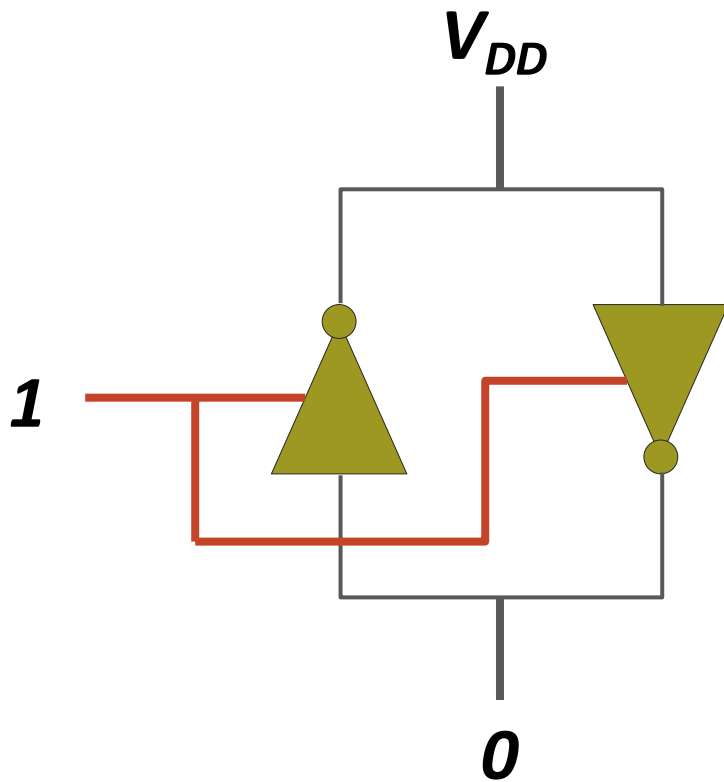
- Energy

- Reliability

- …

# DRAM Chip

# Sense Amplifier



top

enable

bottom

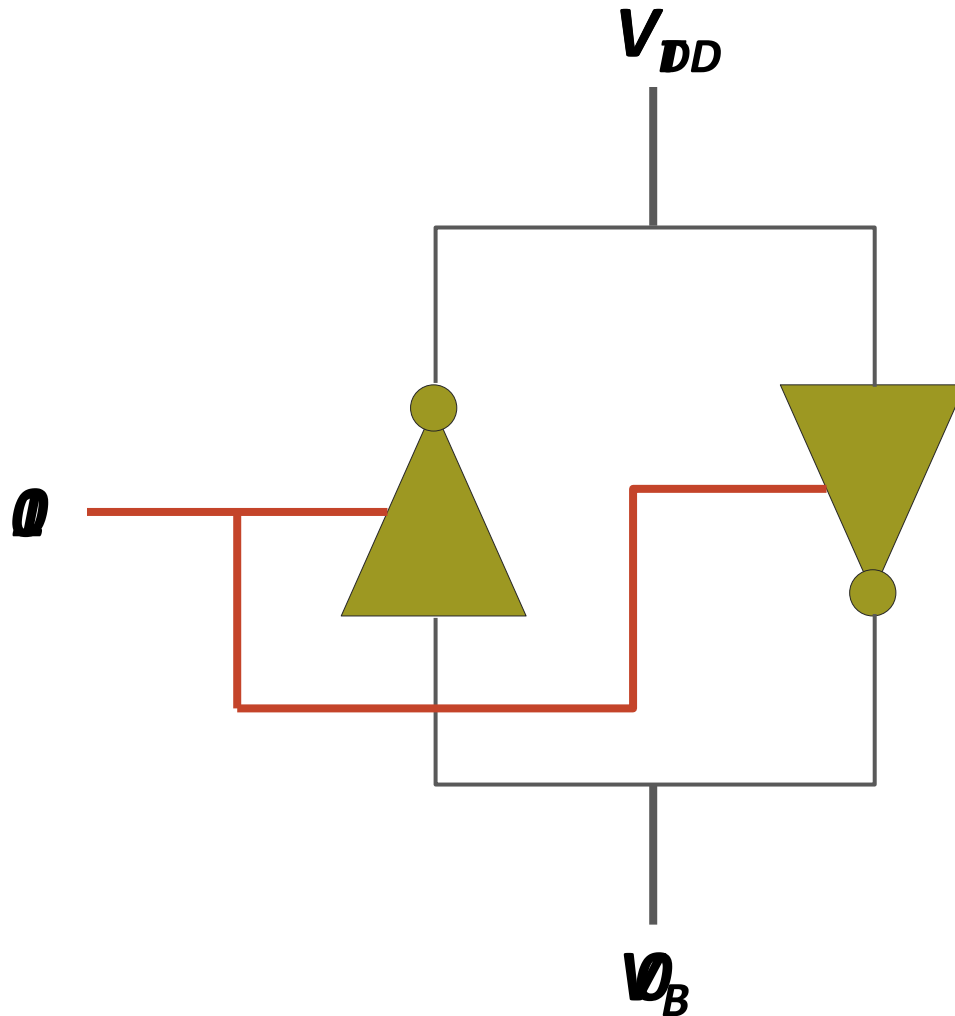Inverter

# Sense Amplifier – Two Stable States



Logical "1"                    Logical "0"

# Sense Amplifier Operation
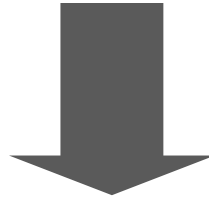


$V_{DD}$

$V_T > V_B$

# DRAM Cell – Capacitor

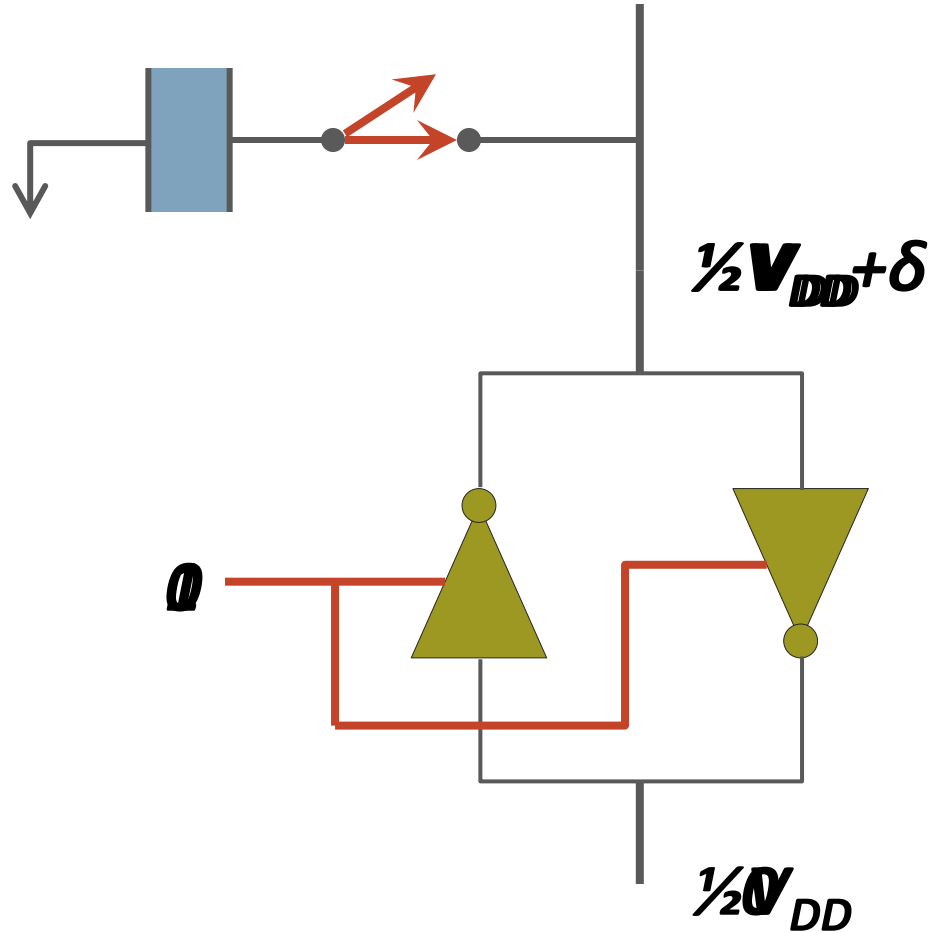Empty State

**Logical "0"**

Fully Charged State

**Logical "1"**

**1** Small – Cannot drive circuits

**2** Reading destroys the state

# Capacitor to Sense Amplifier
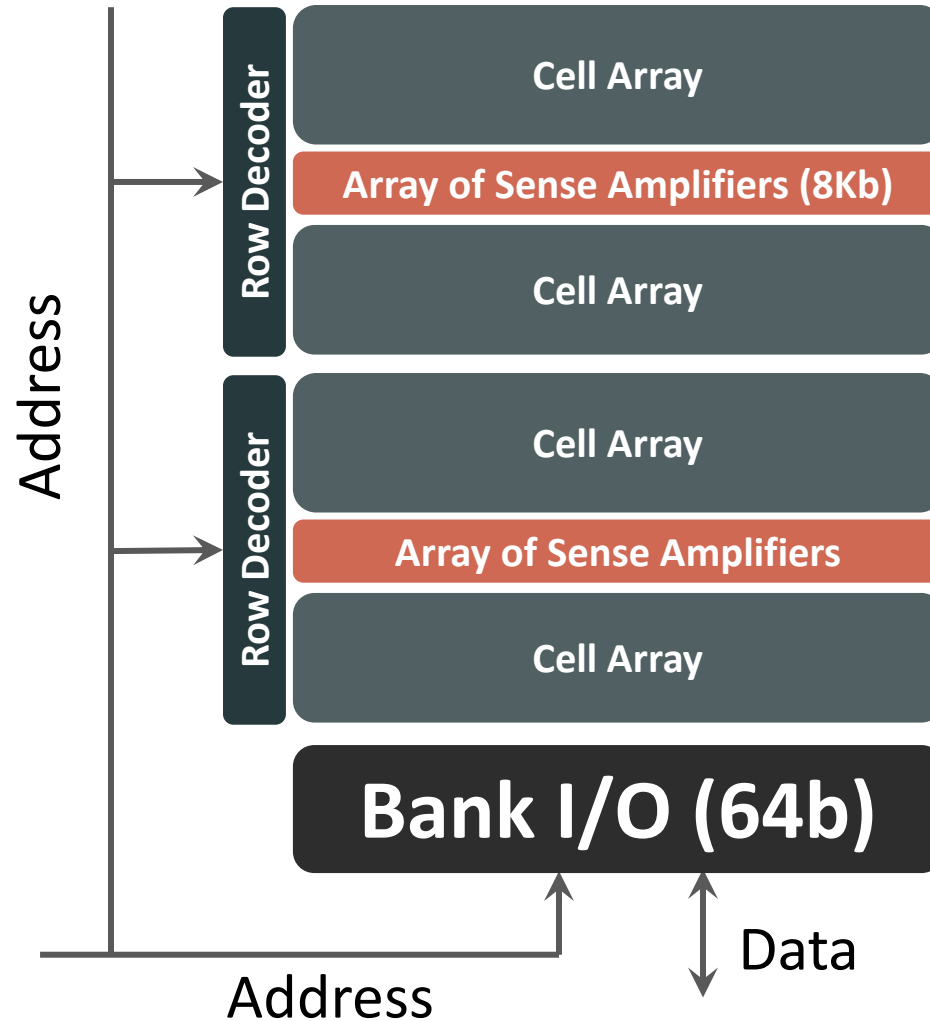
# DRAM Cell Operation



$½V_{DD}+\delta$

$0$

$½V_{DD}$

# DRAM Subarray – Building Block for DRAM Chip

# DRAM Bank

# DRAM Chip

Shared internal bus



Memory channel - 8bits

# DRAM Operation



**1** ACTIVATE Row

**2** READ/WRITE Column

**3** PRECHARGE

Row Address

Row Decoder

Row Decoder

Cell Array

Array of Sense Amplifiers

Cell Array

Bank I/O

Data

Column Address

# Some Key Readings on DRAM (I)

- **DRAM Organization and Operation**

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

    https://people.inf.ethz.ch/omutlu/pub/tldram_hpca13.pdf

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

    https://people.inf.ethz.ch/omutlu/pub/salp-dram_isca12.pdf

  - Lee et al., "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost," ACM TACO 2016.

    https://people.inf.ethz.ch/omutlu/pub/smla_high-bandwidth-3d-stacked-memory_taco16.pdf

# Some Key Readings on DRAM (II)

- **DRAM Refresh**

  - ❑ Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
    https://people.inf.ethz.ch/omutlu/pub/raidr-dram-refresh_isca12.pdf

  - ❑ Chang et al., "Improving DRAM Performance by Parallelizing Refreshes with Accesses," HPCA 2014.
    https://people.inf.ethz.ch/omutlu/pub/dram-access-refresh-parallelization_hpca14.pdf

  - ❑ Patel et al., "The Reach Profiler (REAPER): Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions," ISCA 2017.
    https://people.inf.ethz.ch/omutlu/pub/reaper-dram-retention-profiling-lpddr4_isca17.pdf

# Reading on Simulating Main Memory

- How to evaluate future main memory systems?
- An open-source simulator and its brief description

- Yoongu Kim, Weikun Yang, and Onur Mutlu,
  **"Ramulator: A Fast and Extensible DRAM Simulator"**
  *IEEE Computer Architecture Letters* (**CAL**), March 2015.
  [Source Code]

**SAFARI**

# Memory Systems
## and Memory-Centric Computing Systems

## Lecture 2a: Memory Controllers

Prof. Onur Mutlu

omutlu@gmail.com

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

13 June 2019

TU Wien Fast Course 2019

**SAFARI**　　**ETH** *zürich*　　**Carnegie Mellon**