# Scalable Many-Core Memory Systems Lecture 1, Topic 1: DRAM Basics and DRAM Scaling

Prof. Onur Mutlu
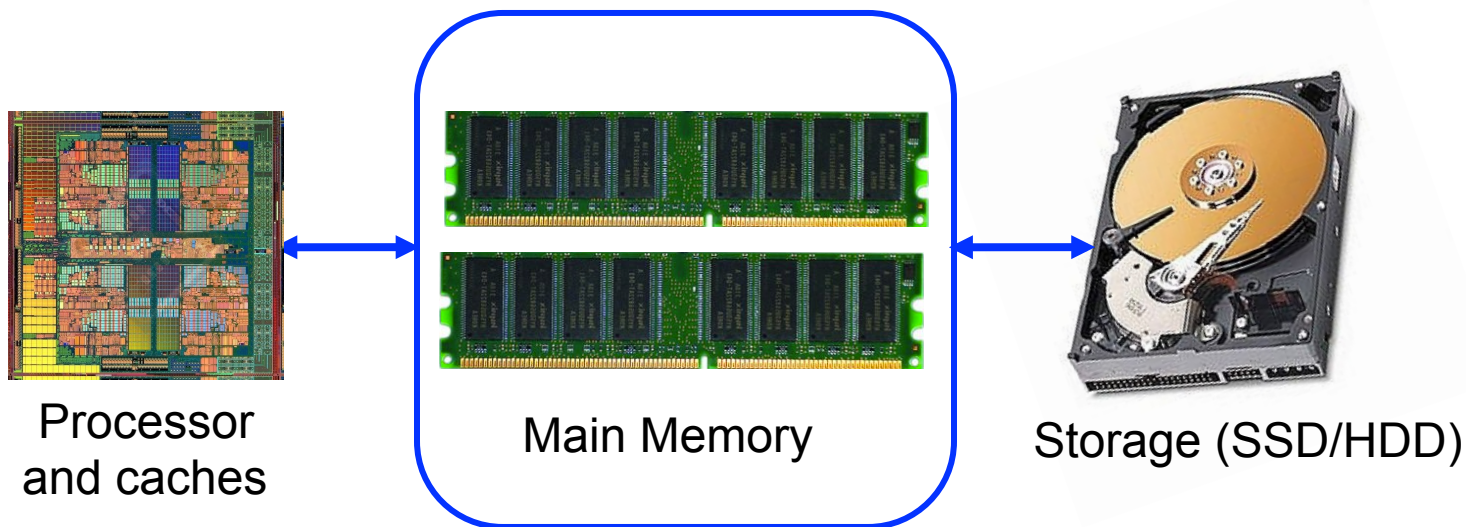
http://www.ece.cmu.edu/~omutlu

onur@cmu.edu

HiPEAC ACACES Summer School 2013

July 15, 2013

**Carnegie Mellon**

**SAFARI**

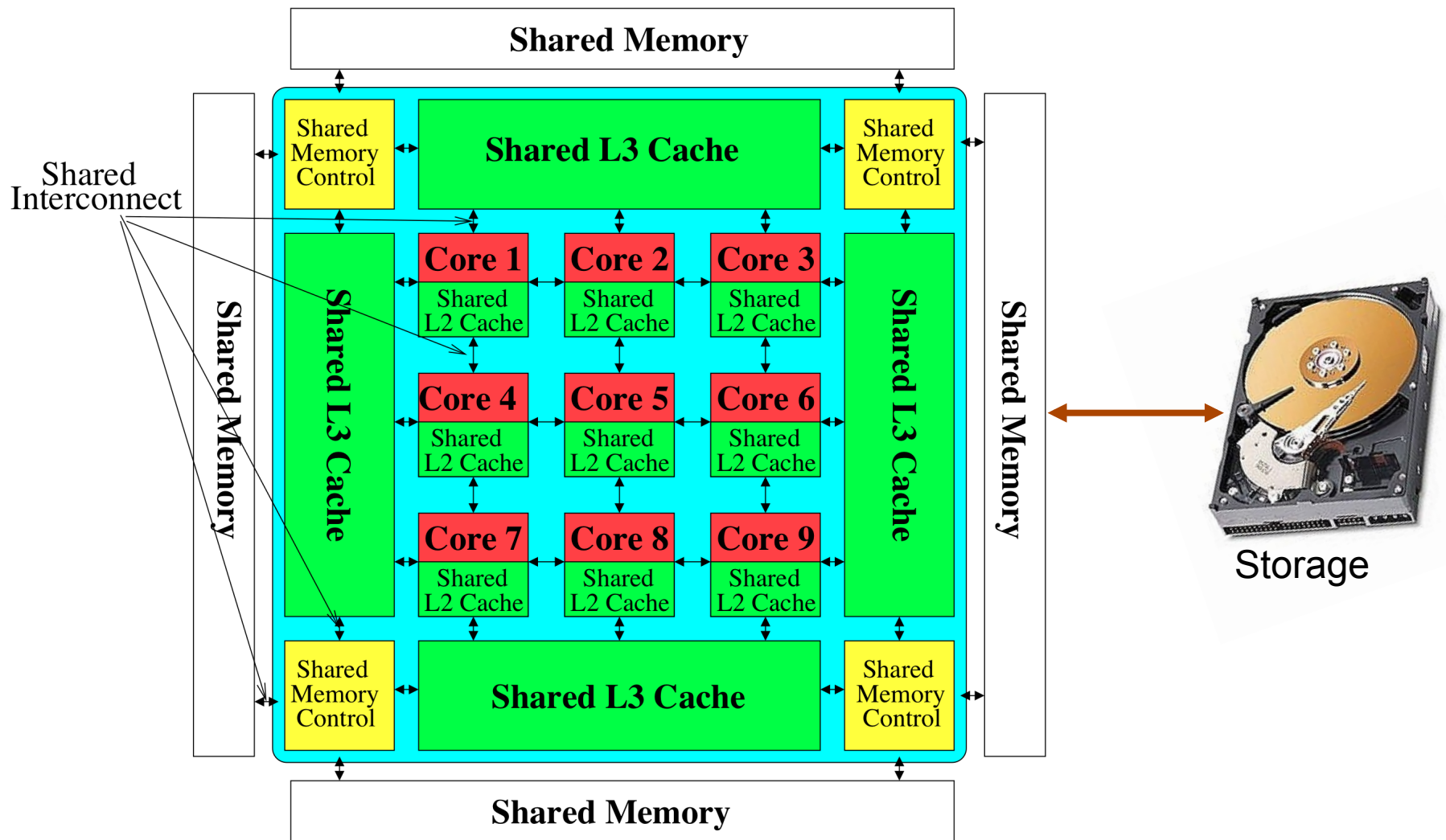# The Main Memory System



Processor and caches — Main Memory — Storage (SSD/HDD)

- **Main memory is a critical component of all computing systems**: server, mobile, embedded, desktop, sensor

- **Main memory system must scale** (in *size*, *technology*, *efficiency*, *cost*, and *management algorithms*) to maintain performance growth and technology scaling benefits

# Memory System: A *Shared Resource* View

# State of the Main Memory System

- Recent technology, architecture, and application trends
  - lead to new requirements
  - exacerbate old requirements

- DRAM and memory controllers, as we know them today, are (will be) unlikely to satisfy all requirements

- Some emerging non-volatile memory technologies (e.g., PCM) enable new opportunities: memory+storage merging

- We need to rethink the main memory system
  - to fix DRAM issues and enable emerging technologies
  - to satisfy all requirements

*SAFARI*

# Major Trends Affecting Main Memory (I)

- Need for main memory capacity, bandwidth, QoS increasing

- Main memory energy/power is a key system design concern

- DRAM technology scaling is ending

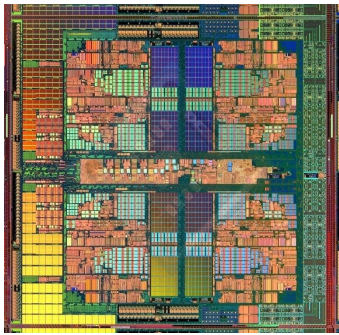**SAFARI**

# Major Trends Affecting Main Memory (II)

- **Need for main memory capacity, bandwidth, QoS increasing**
  - ❑ **Multi-core**: increasing number of cores
  - ❑ **Data-intensive applications**: increasing demand/hunger for data
  - ❑ **Consolidation**: cloud computing, GPUs, mobile
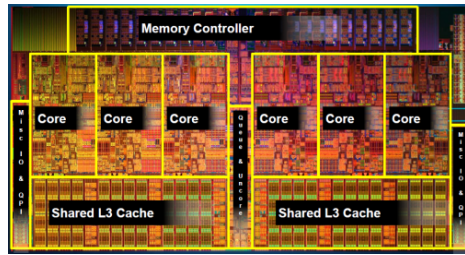
- Main memory energy/power is a key system design concern

- DRAM technology scaling is ending
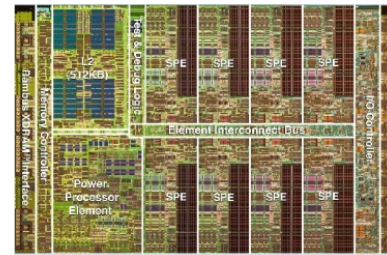
# Example Trend: Many Cores on Chip

- Simpler and lower power than a single large core
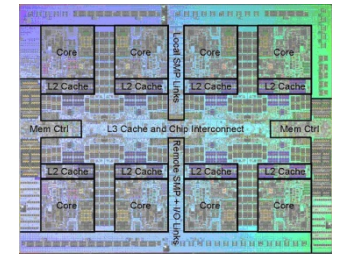- Large scale parallelism on chip

AMD Barcelona
4 cores

Intel Core i7
8 cores

IBM Cell BE
8+1 cores

IBM POWER7
8 cores

Sun Niagara II
8 cores

Nvidia Fermi
448 "cores"

Intel SCC
48 cores, networked

Tilera TILE Gx
100 cores, networked

# Consequence: The Memory Capacity Gap

Core count doubling ~ every 2 years
DRAM DIMM capacity doubling ~ every 3 years



Source: Lim et al., ISCA 2009.

- *Memory capacity per core* expected to drop by 30% every two years
- Trends worse for *memory bandwidth per core*!

# Major Trends Affecting Main Memory (III)

- Need for main memory capacity, bandwidth, QoS increasing

- Main memory energy/power is a key system design concern
  - ~40-50% energy spent in off-chip memory hierarchy [Lefurgy, IEEE Computer 2003]
  - DRAM consumes power even when not used (periodic refresh)

- DRAM technology scaling is ending

**SAFARI**

# Major Trends Affecting Main Memory (IV)

- Need for main memory capacity, bandwidth, QoS increasing

- Main memory energy/power is a key system design concern

- DRAM technology scaling is ending
  - ITRS projects DRAM will not scale easily below X nm
  - Scaling has provided many benefits:
    - higher capacity (density), lower cost, lower energy

**SAFARI**

# The DRAM Scaling Problem

- DRAM stores charge in a capacitor (charge-based memory)
  - Capacitor must be large enough for reliable sensing
  - Access transistor should be large enough for low leakage and high retention time
  - Scaling beyond 40-35nm (2013) is challenging [ITRS, 2009]



- DRAM capacity, cost, and energy/power hard to scale

# Solutions to the DRAM Scaling Problem

- **Two potential solutions**
  - Tolerate DRAM (by taking a fresh look at it)
  - Enable emerging memory technologies to eliminate/minimize DRAM

- **Do both**
  - Hybrid memory systems

# Solution 1: Tolerate DRAM

- Overcome DRAM shortcomings with
  - System-DRAM co-design
  - Novel DRAM architectures, interface, functions
  - Better waste management (efficient utilization)

- Key issues to tackle
  - Reduce refresh energy
  - Improve bandwidth and latency
  - Reduce waste
  - Enable reliability at low cost

- Liu, Jaiyen, Veras, Mutlu, "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.
- Kim, Seshadri, Lee+, "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012.
- Lee+, "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," HPCA 2013.
- Liu+, "An Experimental Study of Data Retention Behavior in Modern DRAM Devices" ISCA'13.
- Seshadri+, "RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," 2013.

# Solution 2: Emerging Memory Technologies

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

- Example: Phase Change Memory
  - Expected to scale to 9nm (2022 [ITRS])
  - Expected to be denser than DRAM: can store multiple bits/cell

- But, emerging technologies have shortcomings as well
  - Can they be enabled to replace/augment/surpass DRAM?

- Lee, Ipek, Mutlu, Burger, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA 2009, CACM 2010, Top Picks 2010.
- Meza, Chang, Yoon, Mutlu, Ranganathan, "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters 2012.
- Yoon, Meza et al., "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012 Best Paper Award.

**SAFARI**

# Hybrid Memory Systems



**CPU**

**DRAM Ctrl**  **PCM Ctrl**

**DRAM**

Fast, **durable**
Small,
leaky, volatile,
high-cost

**Phase Change Memory (or Tech. X)**

Large, non-volatile, low-cost
Slow, **wears out,** high active energy

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

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

# An Orthogonal Issue: Memory Interference

- Problem: Memory interference is uncontrolled → uncontrollable, unpredictable, vulnerable system

- Goal: We need to control it → Design a QoS-aware system

- Solution: Hardware/software cooperative memory QoS
  - Hardware designed to provide a configurable fairness substrate
    - Application-aware memory scheduling, partitioning, throttling
  - Software designed to configure the resources to satisfy different QoS goals

  - E.g., fair, programmable memory controllers and on-chip networks provide QoS and predictable performance
    **[2007-2012, Top Picks'09,'11a,'11b,'12]**

# Agenda for Topic 1 (DRAM Scaling)

- What Will You Learn in This Course
- Main Memory Basics (with a Focus on DRAM)
- Major Trends Affecting Main Memory
- DRAM Scaling Problem and Solution Directions
- Solution Direction 1: System-DRAM Co-Design
- Ongoing Research
- Summary

**SAFARI**

# What Will You Learn in This Course?

- Scalable Many-Core Memory Systems
  - July 15-19, 2013

- Topic 1: Main memory basics, DRAM scaling
- Topic 2: Emerging memory technologies and hybrid memories
- Topic 3: Main memory interference and QoS
- Topic 4 (unlikely): Cache management
- Topic 5 (unlikely): Interconnects

- Major Overview Reading:
  - Mutlu, "Memory Scaling: A Systems Architecture Perspective," IMW 2013.

# This Course

- Will cover many problems and potential solutions related to the design of memory systems in the many core era

- The design of the memory system poses many
  - Difficult research and engineering problems
  - Important fundamental problems
  - Industry-relevant problems

- Many creative and insightful solutions are needed to solve these problems

- Goal: Acquire the basics to develop such solutions (by covering fundamentals and cutting edge research)

**SAFARI**

# Course Information

- My Contact Information
  - Onur Mutlu
  - onur@cmu.edu
  - http://users.ece.cmu.edu/~omutlu
  - +1-512-658-0891 (my cell phone)
  - Find me during breaks and/or email any time.

- Website for Course Slides and Papers
  - http://users.ece.cmu.edu/~omutlu/acaces2013-memory.html
  - http://users.ece.cmu.edu/~omutlu

**SAFARI**

# Readings and Videos

# Overview Reading

- Mutlu, "Memory Scaling: A Systems Architecture Perspective," IMW 2013.


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

# Online Slides (Longer Versions)

- Topic 1: DRAM Basics and DRAM Scaling

  - http://users.ece.cmu.edu/~omutlu/pub/onur-ACACES2013-Topic1-dram-basics-and-scaling.pptx

  - http://users.ece.cmu.edu/~omutlu/pub/onur-ACACES2013-Topic1-dram-basics-and-scaling.pdf

- Topic 2: Emerging Technologies and Hybrid Memories

  - http://users.ece.cmu.edu/~omutlu/pub/onur-ACACES2013-Topic2-emerging-and-hybrid-memory-technologies.pptx

  - http://users.ece.cmu.edu/~omutlu/pub/onur-ACACES2013-Topic2-emerging-and-hybrid-memory-technologies.pdf

- Topic 3: Memory Interference and QoS-Aware Memory Systems

  - http://users.ece.cmu.edu/~omutlu/pub/onur-ACACES2013-Topic3-memory-qos.pptx

  - http://users.ece.cmu.edu/~omutlu/pub/onur-ACACES2013-Topic3-memory-qos.pdf

# Memory Lecture Videos

- Memory Hierarchy (and Introduction to Caches)
  - http://www.youtube.com/watch?v=JBdfZ5i21cs&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=22
- Main Memory
  - http://www.youtube.com/watch?v=ZLCy3pG7Rc0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=25
- Memory Controllers, Memory Scheduling, Memory QoS
  - http://www.youtube.com/watch?v=ZSotvL3WXmA&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=26
  - http://www.youtube.com/watch?v=1xe2w3_NzmI&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=27
- Emerging Memory Technologies
  - http://www.youtube.com/watch?v=LzfOghMKyA0&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=35
- Multiprocessor Correctness and Cache Coherence
  - http://www.youtube.com/watch?v=U-VZKMgItDM&list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ&index=32

# Readings for Topic 1 (DRAM Scaling)

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

- Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," ISCA 2012.

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

- Liu et al., "An Experimental Study of Data Retention Behavior in Modern DRAM Devices," ISCA 2013.

- Seshadri et al., "RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data," CMU CS Tech Report 2013.

- David et al., "Memory Power Management via Dynamic Voltage/Frequency Scaling," ICAC 2011.

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

# Readings for Topic 2 (Emerging Technologies)

- Lee, Ipek, Mutlu, Burger, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA 2009, CACM 2010, Top Picks 2010.

- Qureshi et al., "Scalable high performance main memory system using phase-change memory technology," ISCA 2009.

- Meza et al., "Enabling Efficient and Scalable Hybrid Memories," IEEE Comp. Arch. Letters 2012.

- Yoon et al., "Row Buffer Locality Aware Caching Policies for Hybrid Memories," ICCD 2012 Best Paper Award.

- Meza et al., "A Case for Efficient Hardware-Software Cooperative Management of Storage and Memory," WEED 2013.

- Kultursay et al., "Evaluating STT-RAM as an Energy-Efficient Main Memory Alternative," ISPASS 2013.

- Cai et al., "Error Analysis and Retention-Aware Error Management for NAND Flash Memory," ITJ 2013.

# Readings for Topic 3 (Memory QoS)

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

- Mutlu and Moscibroda, "Stall-Time Fair Memory Access Scheduling," MICRO 2007.

- Mutlu and Moscibroda, "Parallelism-Aware Batch Scheduling," ISCA 2008, IEEE Micro 2009.

- Kim et al., "ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers," HPCA 2010.

- Kim et al., "Thread Cluster Memory Scheduling," MICRO 2010, IEEE Micro 2011.

- Muralidhara et al., "Memory Channel Partitioning," MICRO 2011.

- Ausavarungnirun et al., "Staged Memory Scheduling," ISCA 2012.

- Subramanian et al., "MISE: Providing Performance Predictability and Improving Fairness in Shared Main Memory Systems," HPCA 2013.

- Das et al., "Application-to-Core Mapping Policies to Reduce Memory System Interference in Multi-Core Systems," HPCA 2013.

# Readings for Topic 3 (Memory QoS)

- Ebrahimi et al., "Fairness via Source Throttling," ASPLOS 2010, ACM TOCS 2012.

- Lee et al., "Prefetch-Aware DRAM Controllers," MICRO 2008, IEEE TC 2011.

- Ebrahimi et al., "Parallel Application Memory Scheduling," MICRO 2011.

- Ebrahimi et al., "Prefetch-Aware Shared Resource Management for Multi-Core Systems," ISCA 2011.

# Readings in Flash Memory

- Yu Cai, Gulay Yalcin, <u>Onur Mutlu</u>, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai,
  **"Error Analysis and Retention-Aware Error Management for NAND Flash Memory"**
  *Intel Technology Journal* (**ITJ**) *Special Issue on Memory Resiliency*, Vol. 17, No. 1, May 2013.

- Yu Cai, Erich F. Haratsch, <u>Onur Mutlu</u>, and Ken Mai,
  **"Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis and Modeling"**
  *Proceedings of the Design, Automation, and Test in Europe Conference* (**DATE**), Grenoble, France, March 2013. Slides (ppt)

- Yu Cai, Gulay Yalcin, <u>Onur Mutlu</u>, Erich F. Haratsch, Adrian Cristal, Osman Unsal, and Ken Mai,
  **"Flash Correct-and-Refresh: Retention-Aware Error Management for Increased Flash Memory Lifetime"**
  *Proceedings of the 30th IEEE International Conference on Computer Design* (**ICCD**), Montreal, Quebec, Canada, September 2012. Slides (ppt) (pdf)

- Yu Cai, Erich F. Haratsch, <u>Onur Mutlu</u>, and Ken Mai,
  **"Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis"**
  *Proceedings of the Design, Automation, and Test in Europe Conference* (**DATE**), Dresden, Germany, March 2012. Slides (ppt)

# Online Lectures and More Information

- Online Computer Architecture Lectures
  - http://www.youtube.com/playlist?list=PL5PHm2jkkXmidJOd59REog9jDnPDTG6IJ

- Online Computer Architecture Courses
  - Intro: http://www.ece.cmu.edu/~ece447/s13/doku.php
  - Advanced: http://www.ece.cmu.edu/~ece740/f11/doku.php
  - Advanced: http://www.ece.cmu.edu/~ece742/doku.php

- Recent Research Papers
  - http://users.ece.cmu.edu/~omutlu/projects.htm
  - http://scholar.google.com/citations?user=7XyGUGkAAAAJ&hl=en

# Agenda for Topic 1 (DRAM Scaling)

- What Will You Learn in This Mini-Lecture Series
- Main Memory Basics (with a Focus on DRAM)
- Major Trends Affecting Main Memory
- DRAM Scaling Problem and Solution Directions
- Solution Direction 1: System-DRAM Co-Design
- Ongoing Research
- Summary

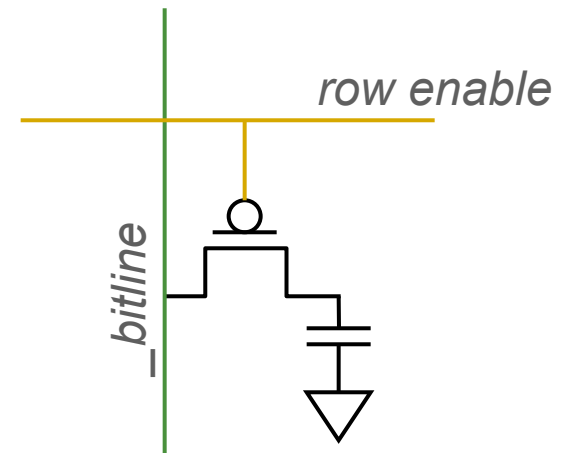# Main Memory

# Main Memory in the System

# Ideal Memory

- Zero access time (latency)
- Infinite capacity
- Zero cost
- Infinite bandwidth (to support multiple accesses in parallel)

# The Problem

- Ideal memory's requirements oppose each other

- Bigger is slower
  - Bigger $\rightarrow$ Takes longer to determine the location

- Faster is more expensive
  - Memory technology: SRAM vs. DRAM

- Higher bandwidth is more expensive
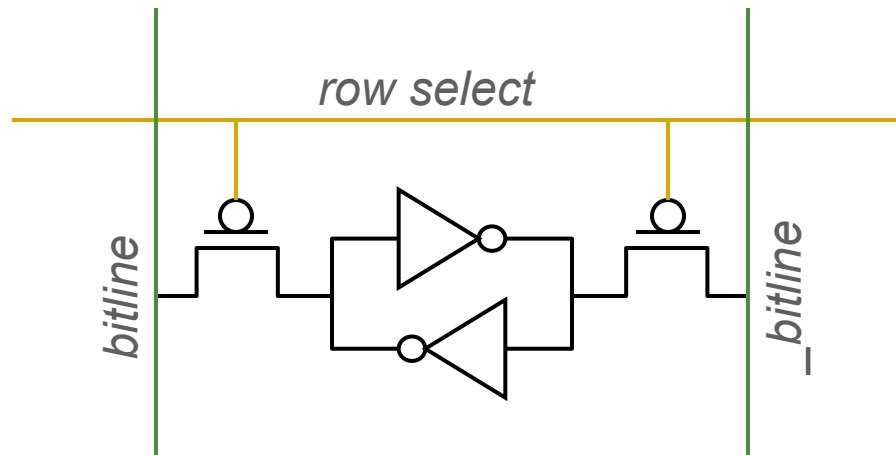  - Need more banks, more ports, higher frequency, or faster technology

# Memory Technology: DRAM

- Dynamic random access memory

- Capacitor charge state indicates stored value
  - Whether the capacitor is charged or discharged indicates storage of 1 or 0
  - 1 capacitor
  - 1 access transistor

  *row enable*

  *_bitline*

- Capacitor leaks through the RC path
  - DRAM cell loses charge over time
  - DRAM cell needs to be refreshed

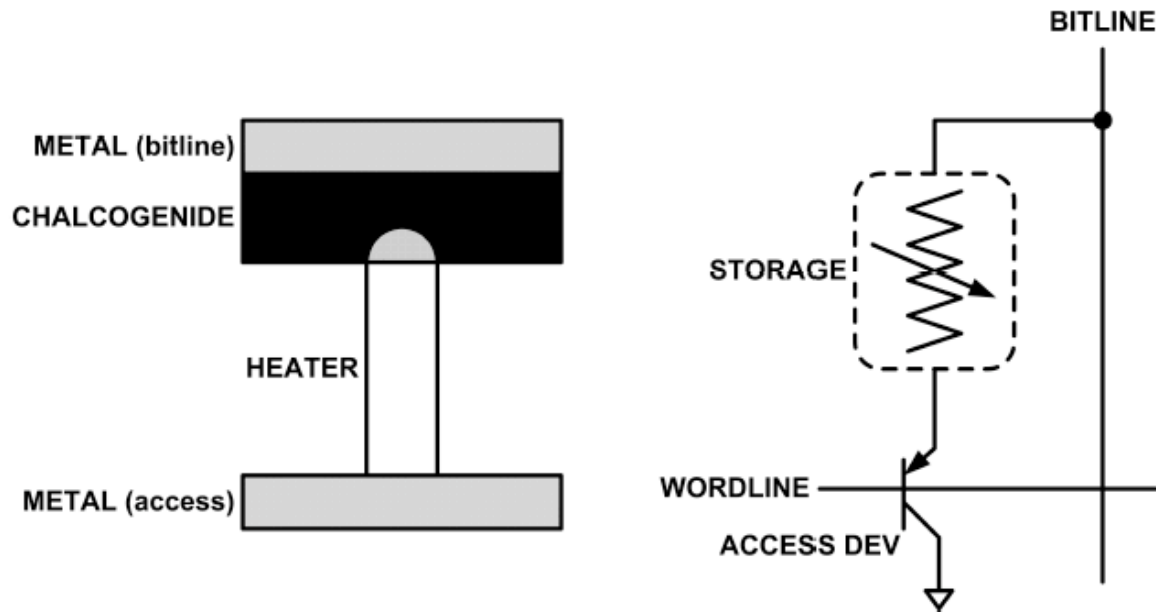  - Read Liu et al., "RAIDR: Retention-aware Intelligent DRAM Refresh," ISCA 2012.

# Memory Technology: SRAM

- Static random access memory
- Two cross coupled inverters store a single bit
    - Feedback path enables the stored value to persist in the "cell"
    - 4 transistors for storage
    - 2 transistors for access

*row select*

*bitline*

*_bitline*

# An Aside: Phase Change Memory

- Phase change material (chalcogenide glass) exists in two states:
  - Amorphous: Low optical reflexivity and high electrical resistivity
  - Crystalline: High optical reflexivity and low electrical resistivity

BITLINE

METAL (bitline)

CHALCOGENIDE

HEATER

METAL (access)

STORAGE

WORDLINE

ACCESS DEV

PCM is resistive memory:  High resistance (0), Low resistance (1)

Lee, Ipek, Mutlu, Burger, "Architecting Phase Change Memory as a Scalable DRAM Alternative," ISCA 2009.

# Memory Bank: A Fundamental Concept

- Interleaving (banking)
  - Problem: a single monolithic memory array takes long to access and does not enable multiple accesses in parallel

  - Goal: Reduce the latency of memory array access and enable multiple accesses in parallel

  - Idea: Divide the array into multiple banks that can be accessed independently (in the same cycle or in consecutive cycles)
    - Each bank is smaller than the entire memory storage
    - Accesses to different banks can be overlapped

  - An issue: How do you map data to different banks? (i.e., how do you interleave data across banks?)

# Memory Bank Organization and Operation



- Read access sequence:

  1. Decode row address & drive word-lines

  2. Selected bits drive bit-lines
     - Entire row read

  3. Amplify row data

  4. Decode column address & select subset of row
     - Send to output

  5. Precharge bit-lines
     - For next access

# Why Memory Hierarchy?

- We want both fast and large

- But we cannot achieve both with a single level of memory

- Idea: Have multiple levels of storage (progressively bigger and slower as the levels are farther from the processor) and ensure most of the data the processor needs is kept in the fast(er) level(s)

# Memory Hierarchy

- **Fundamental tradeoff**
  - Fast memory: small
  - Large memory: slow
- Idea: Memory hierarchy

```
+--------+  +--------+      +-----------+   +-----------+
|  CPU   |  |        |<---->|   Main    |   |           |
| +----+ |<>| Cache  |      |  Memory   |   | Hard Disk |
| | RF | |  |        |      |  (DRAM)   |   |           |
| +----+ |  |        |      |           |   |           |
+--------+  +--------+      +-----------+   +-----------+
```

- **Latency, cost, size, bandwidth**

# Caching Basics: Exploit Temporal Locality

- Idea: Store recently accessed data in automatically managed fast memory (called cache)
- Anticipation: the data will be accessed again soon

- Temporal locality principle
  - Recently accessed data will be again accessed in the near future
  - This is what Maurice Wilkes had in mind:
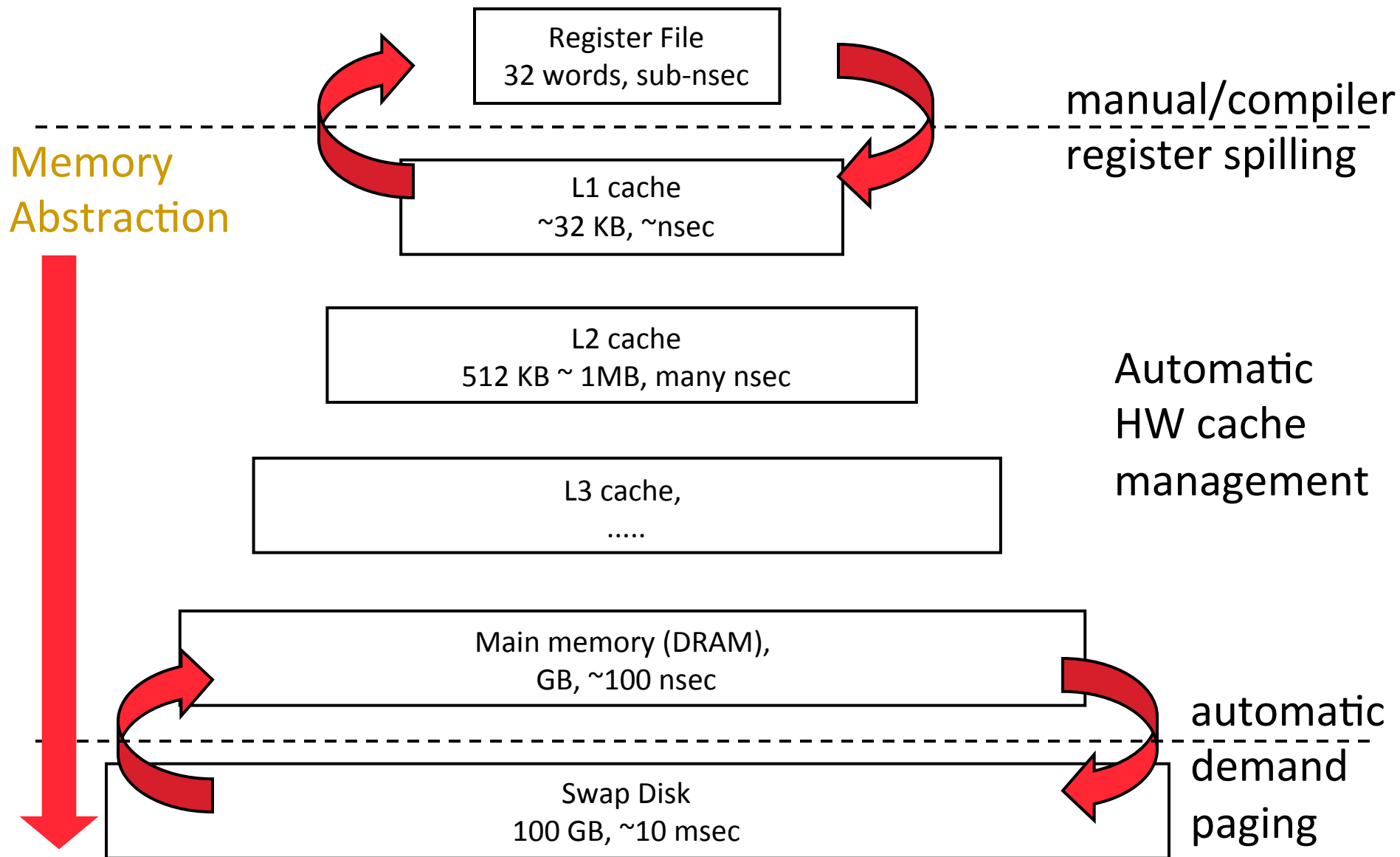    - Wilkes, "Slave Memories and Dynamic Storage Allocation," IEEE Trans. On Electronic Computers, 1965.
    - "The use is discussed of a fast core memory of, say 32000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory."

# Caching Basics: Exploit Spatial Locality

- Idea: Store addresses adjacent to the recently accessed one in automatically managed fast memory
  - Logically divide memory into equal size blocks
  - Fetch to cache the accessed block in its entirety
- Anticipation: nearby data will be accessed soon

- Spatial locality principle
  - Nearby data in memory will be accessed in the near future
    - E.g., sequential instruction access, array traversal
  - This is what IBM 360/85 implemented
    - 16 Kbyte cache with 64 byte blocks
    - Liptay, "Structural aspects of the System/360 Model 85 II: the cache," IBM Systems Journal, 1968.

# A Note on Manual vs. Automatic Management

- **Manual:** Programmer manages data movement across levels
  -- too painful for programmers on substantial programs
  - "core" vs "drum" memory in the 50's
  - still done in some embedded processors (on-chip scratch pad SRAM in lieu of a cache)

- **Automatic:** Hardware manages data movement across levels, transparently to the programmer
  ++ programmer's life is easier
  - simple heuristic: keep most recently used items in cache
  - the average programmer doesn't need to know about it
    - You don't need to know how big the cache is and how it works to write a "correct" program! (What if you want a "fast" program?)

# Automatic Management in Memory Hierarchy

- Wilkes, "Slave Memories and Dynamic Storage Allocation," IEEE Trans. On Electronic Computers, 1965.

## Slave Memories and Dynamic Storage Allocation

### M. V. WILKES

#### SUMMARY

The use is discussed of a fast core memory of, say, 32 000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory.

- "By a slave memory I mean one which automatically accumulates to itself words that come from a slower main memory, and keeps them available for subsequent use without it being necessary for the penalty of main memory access to be incurred again."

# A Modern Memory Hierarchy

Memory
Abstraction

Register File
32 words, sub-nsec

manual/compiler
register spilling

L1 cache
~32 KB, ~nsec

L2 cache
512 KB ~ 1MB, many nsec

Automatic
HW cache
management

L3 cache,
.....

Main memory (DRAM),
GB, ~100 nsec

automatic
demand
paging

Swap Disk
100 GB, ~10 msec

47

# The DRAM Subsystem

# DRAM Subsystem Organization

- Channel
- DIMM
- Rank
- Chip
- Bank
- Row/Column

# Page Mode DRAM

- A DRAM bank is a 2D array of cells: rows x columns
- A "DRAM row" is also called a "DRAM page"
- "Sense amplifiers" also called "row buffer"

- Each address is a <row,column> pair
- Access to a "closed row"
  - Activate command opens row (placed into row buffer)
  - Read/write command reads/writes column in the row buffer
  - Precharge command closes the row and prepares the bank for next access
- Access to an "open row"
  - No need for activate command

# DRAM Bank Operation

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

Columns

Rows

Row decoder

Row address 0 1

Row 1    Row Buffer  CONFLICT !  HIT

Column address 0 85

Column mux

Data

# The DRAM Chip

- Consists of multiple banks (2-16 in Synchronous DRAM)
- Banks share command/address/data buses
- The chip itself has a narrow interface (4-16 bits per read)

# 128M x 8-bit DRAM Chip

# DRAM Rank and Module

- Rank: Multiple chips operated together to form a wide interface

- All chips comprising a rank are controlled at the same time
  - Respond to a single command
  - Share address and command buses, but provide different data

- A DRAM module consists of one or more ranks
  - E.g., DIMM (dual inline memory module)
  - This is what you plug into your motherboard

- If we have chips with 8-bit interface, to read 8 bytes in a single access, use 8 chips in a DIMM

# A 64-bit Wide DIMM (One Rank)

# A 64-bit Wide DIMM (One Rank)



- **Advantages:**
  - Acts like a high-capacity DRAM chip with a wide interface
  - Flexibility: memory controller does not need to deal with individual chips

- **Disadvantages:**
  - Granularity: Accesses cannot be smaller than the interface width

# Multiple DIMMs



- **Advantages:**
  - Enables even higher capacity

- **Disadvantages:**
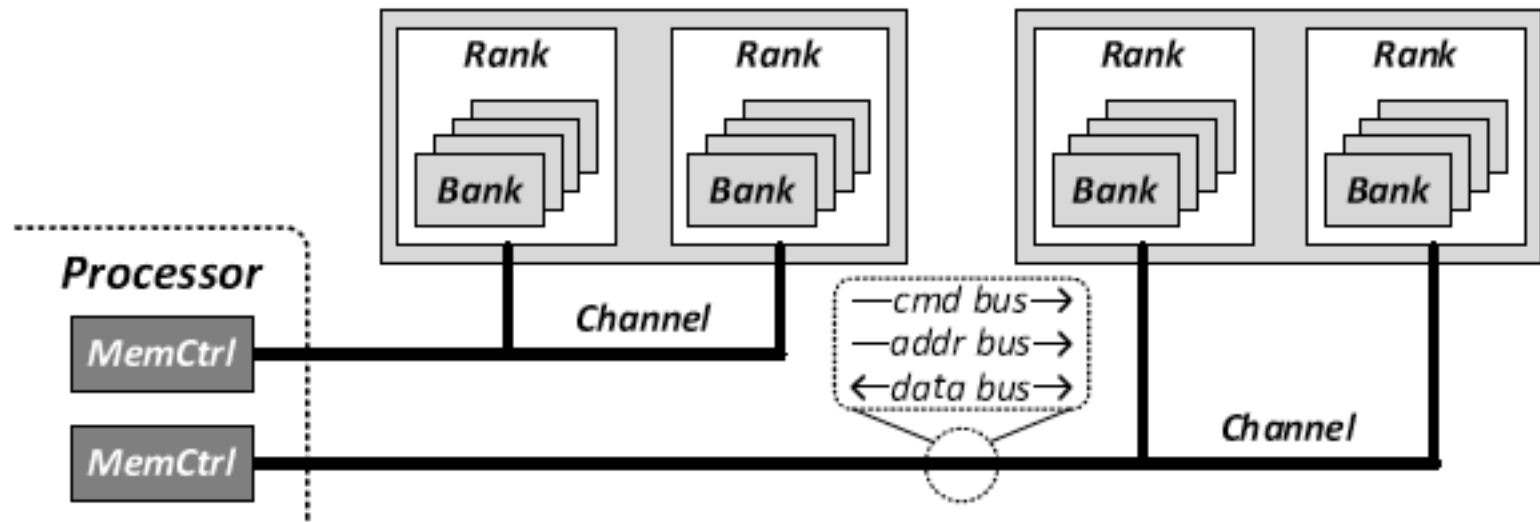  - Interconnect complexity and energy consumption can be high

"Mesh Topology"

Addr & Cmd
Data Bus
Chip (DIMM) Select

# DRAM Channels



- **2 Independent Channels: 2 Memory Controllers (Above)**
- **2 Dependent/Lockstep Channels: 1 Memory Controller with wide interface (Not shown above)**

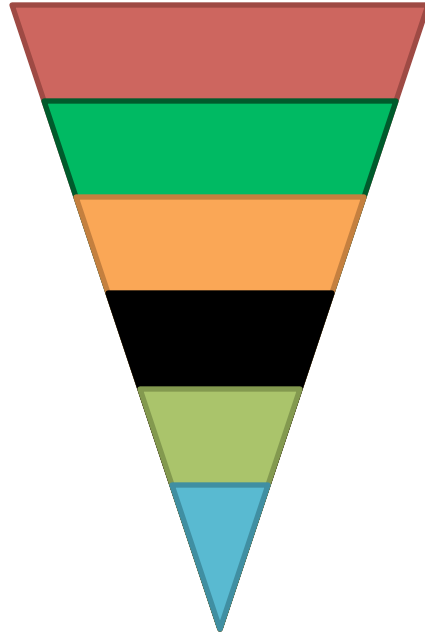# Generalized Memory Structure

# Generalized Memory Structure



Kim+, "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012.

# The DRAM Subsystem
# The Top Down View

# DRAM Subsystem Organization
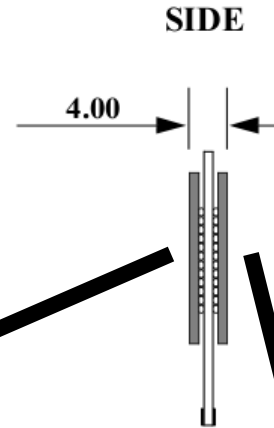
- Channel
- DIMM
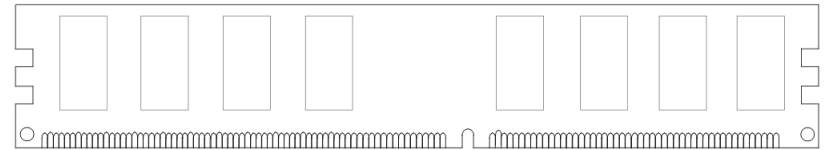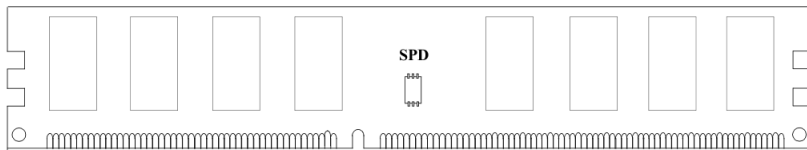- Rank
- Chip
- Bank
- Row/Column

# The DRAM subsystem

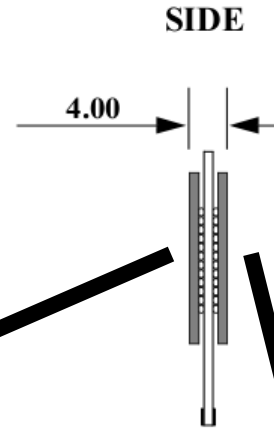# Breaking down a DIMM

**DIMM** **(Dual in-line memory module)**

**SIDE**

4.00

Side view

**Front of DIMM**

SPD

**Back of DIMM**

# Breaking down a DIMM

**DIMM** **(Dual in-line memory module)**

**SIDE**

4.00

Side view

**Front of DIMM**

**Back of DIMM**

SPD

**Rank 0:** collection of 8 chips

**Rank 1**

# Rank



Rank 0 (Front)

Rank 1 (Back)

<0:63>

<0:63>

Addr/Cmd

CS <0:1>

Data <0:63>

**Memory channel**

# Breaking down a Rank

# Breaking down a Chip

Chip 0

8 banks

<0:7>

<0:7>

<0:7>

<0:7>

<0:7>

. . .

# Breaking down a Bank

# DRAM Subsystem Organization

- Channel
- DIMM
- Rank
- Chip
- Bank
- Row/Column

# Example: Transferring a cache block

**Physical memory space**

0xFFFF...F

0x40

0x00

64B
cache block

Mapped to

**Channel 0**

**DIMM 0**

**Rank 0**

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block



A 64B cache block takes 8 I/O cycles to transfer.

During the process, 8 columns are read sequentially.
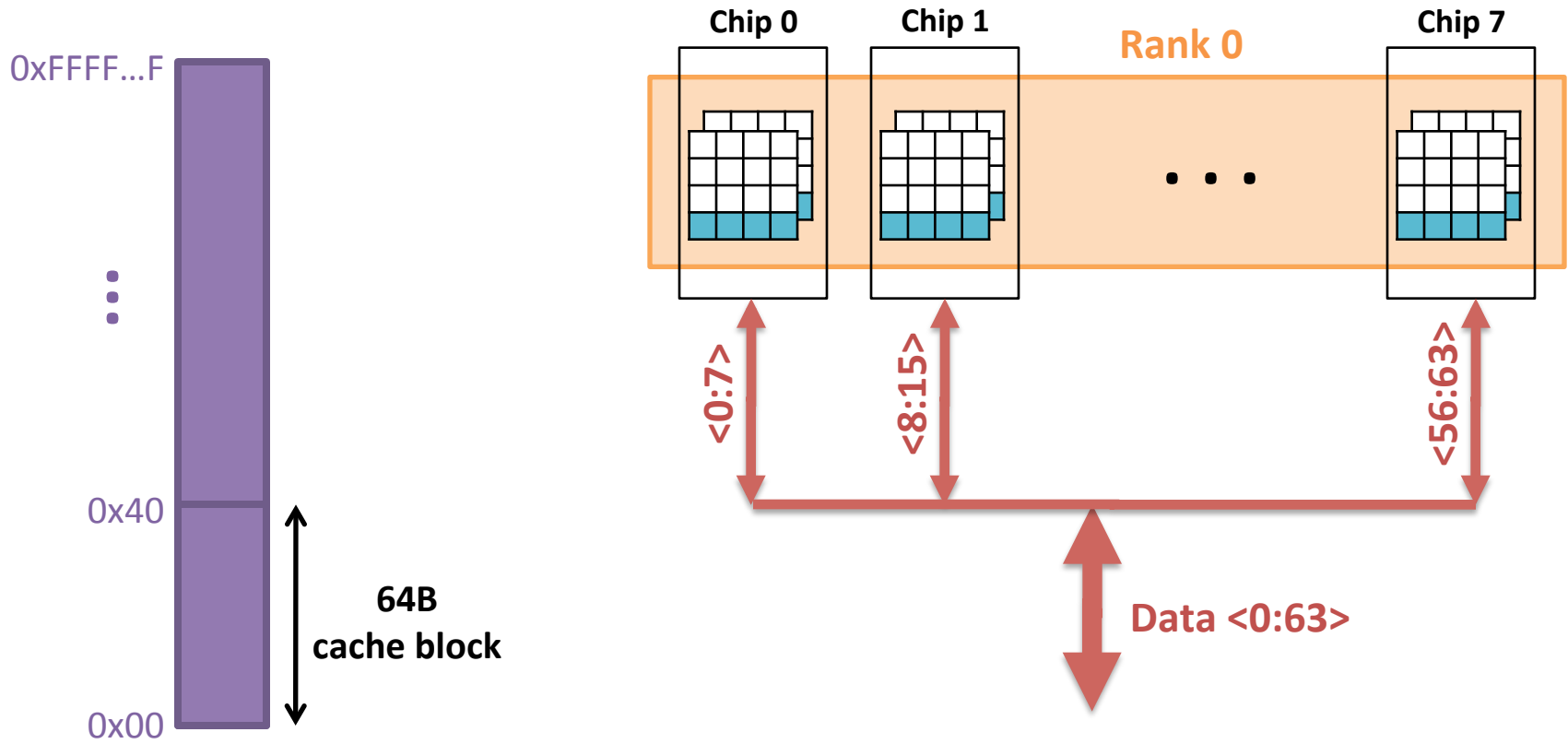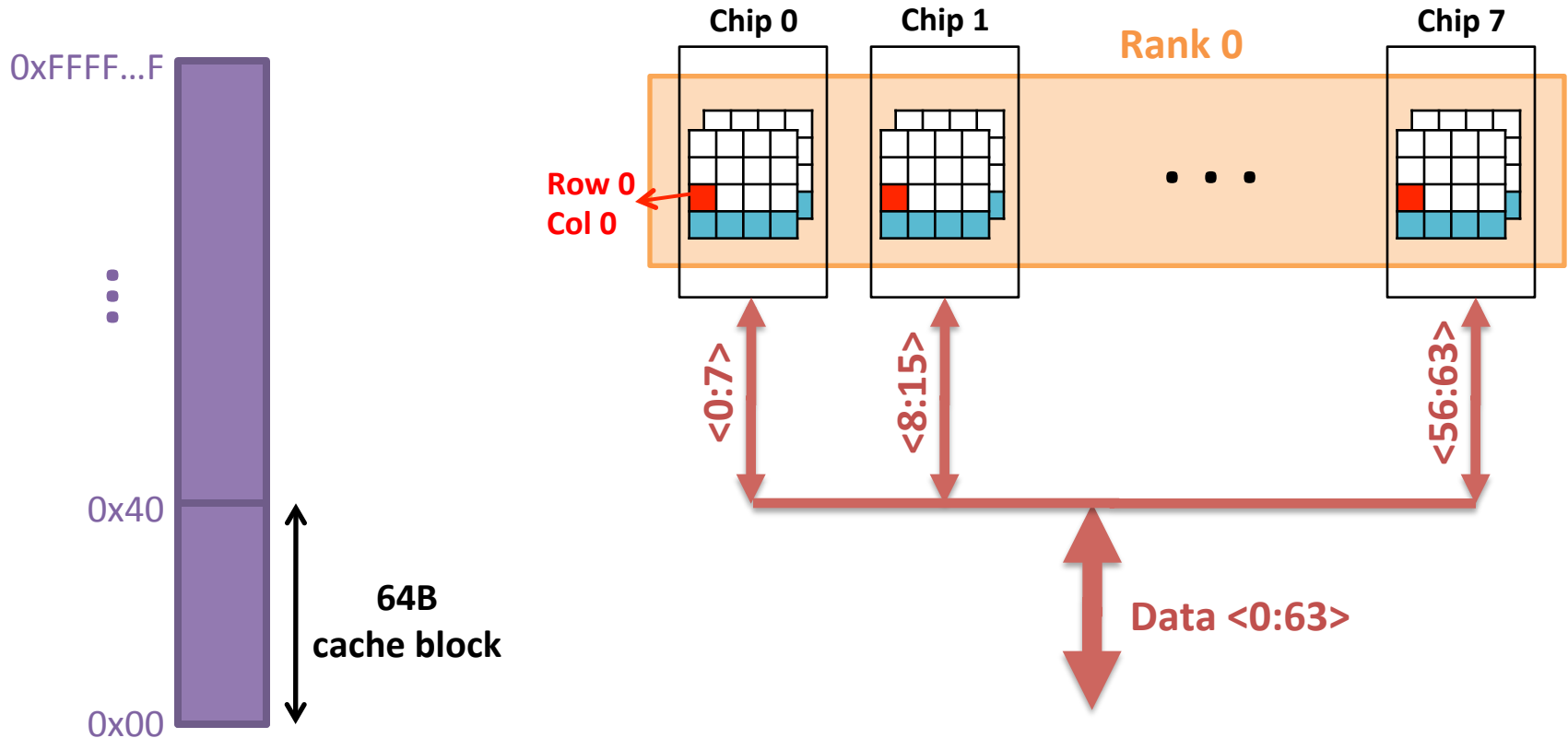
# Latency Components: Basic DRAM Operation

- CPU → controller transfer time
- Controller latency
  - Queuing & scheduling delay at the controller
  - Access converted to basic commands
- Controller → DRAM transfer time
- DRAM bank latency
  - Simple CAS (column address strobe) if row is "open" OR
  - RAS (row address strobe) + CAS if array precharged OR
  - PRE + RAS + CAS (worst case)
- DRAM → Controller transfer time
  - Bus latency (BL)
- Controller to CPU transfer time

# Multiple Banks (Interleaving) and Channels

- Multiple banks
  - Enable concurrent DRAM accesses
  - Bits in address determine which bank an address resides in
- Multiple independent channels serve the same purpose
  - But they are even better because they have separate data buses
  - Increased bus bandwidth

- Enabling more concurrency requires reducing
  - Bank conflicts
  - Channel conflicts
- How to select/randomize bank/channel indices in address?
  - Lower order bits have more entropy
  - Randomizing hash functions (XOR of different address bits)

# How Multiple Banks Help



Before: No Overlapping
Assuming accesses to different DRAM rows

After: Overlapped Accesses
Assuming no bank conflicts

# Address Mapping (Single Channel)

- Single-channel system with 8-byte memory bus
  - 2GB memory, 8 banks, 16K rows & 2K columns per bank

- Row interleaving
  - Consecutive rows of memory in consecutive banks

| Row (14 bits) | Bank (3 bits) | Column (11 bits) | Byte in bus (3 bits) |
|---|---|---|---|

  - Accesses to consecutive cache blocks serviced in a pipelined manner

- Cache block interleaving
  - Consecutive cache block addresses in consecutive banks
  - 64 byte cache blocks

| Row (14 bits) | High Column | Bank (3 bits) | Low Col. | Byte in bus (3 bits) |
|---|---|---|---|---|
| | 8 bits | | 3 bits | |

  - Accesses to consecutive cache blocks can be serviced in parallel

# Bank Mapping Randomization

- DRAM controller can randomize the address mapping to banks so that bank conflicts are less likely

| | 3 bits | Column (11 bits) | Byte in bus (3 bits) |
|---|---|---|---|

XOR

Bank index
(3 bits)

# Address Mapping (Multiple Channels)

| C | Row (14 bits) | Bank (3 bits) | Column (11 bits) | Byte in bus (3 bits) |
|---|---|---|---|---|

| Row (14 bits) | C | Bank (3 bits) | Column (11 bits) | Byte in bus (3 bits) |
|---|---|---|---|---|

| Row (14 bits) | Bank (3 bits) | C | Column (11 bits) | Byte in bus (3 bits) |
|---|---|---|---|---|

| Row (14 bits) | Bank (3 bits) | Column (11 bits) | C | Byte in bus (3 bits) |
|---|---|---|---|---|

- ## Where are consecutive cache blocks?

| C | Row (14 bits) | High Column | Bank (3 bits) | Low Col. | Byte in bus (3 bits) |
|---|---|---|---|---|---|
| | | 8 bits | | 3 bits | |

| Row (14 bits) | C | High Column | Bank (3 bits) | Low Col. | Byte in bus (3 bits) |
|---|---|---|---|---|---|
| | | 8 bits | | 3 bits | |

| Row (14 bits) | High Column | C | Bank (3 bits) | Low Col. | Byte in bus (3 bits) |
|---|---|---|---|---|---|
| | 8 bits | | | 3 bits | |

| Row (14 bits) | High Column | Bank (3 bits) | C | Low Col. | Byte in bus (3 bits) |
|---|---|---|---|---|---|
| | 8 bits | | | 3 bits | |

| Row (14 bits) | High Column | Bank (3 bits) | Low Col. | C | Byte in bus (3 bits) |
|---|---|---|---|---|---|
| | 8 bits | | 3 bits | | |

# Interaction with Virtual➔Physical Mapping

- **Operating System influences where an address maps to in DRAM**

| Virtual Page number (52 bits) | | Page offset (12 bits) | **VA** |
|---|---|---|---|

| Physical Frame number (19 bits) | | Page offset (12 bits) | **PA** |
|---|---|---|---|

| Row (14 bits) | Bank (3 bits) | Column (11 bits) | Byte in bus (3 bits) | **PA** |
|---|---|---|---|---|

- **Operating system can influence which bank/channel/rank a virtual page is mapped to.**

- It can perform page coloring to
  - Minimize bank conflicts
  - Minimize inter-application interference **[Muralidhara+ MICRO'11]**

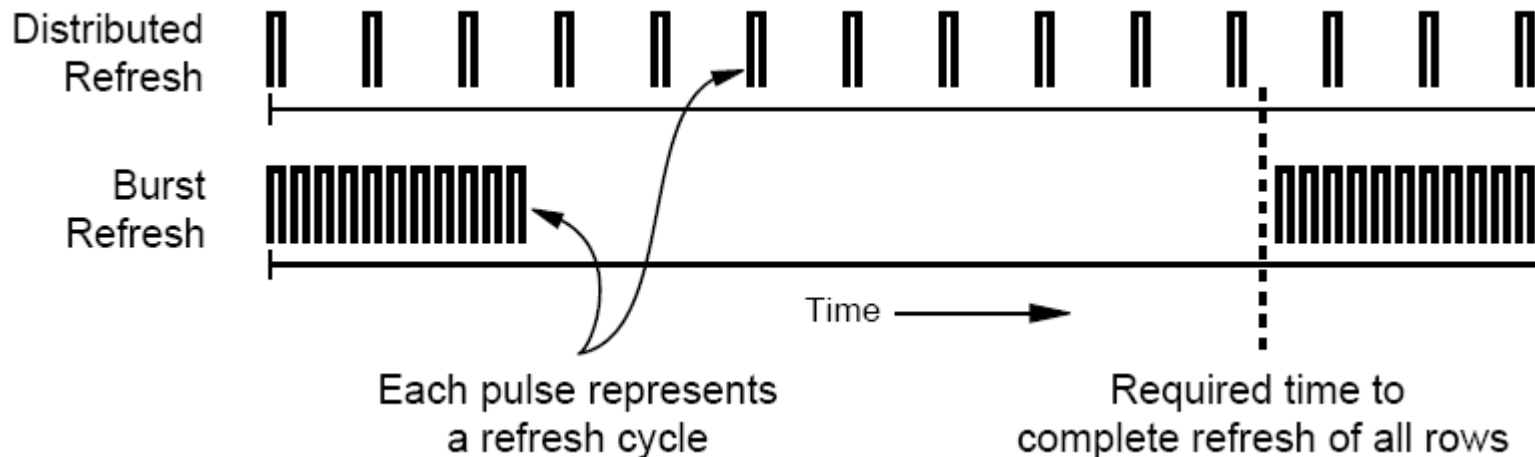# DRAM Refresh (I)

- DRAM capacitor charge leaks over time
- The memory controller needs to read each row periodically to restore the charge
  - ❑ Activate + precharge each row every N ms
  - ❑ Typical N = 64 ms
- Implications on performance?
- -- DRAM bank unavailable while refreshed
- -- Long pause times: If we refresh all rows in burst, every 64ms the DRAM will be unavailable until refresh ends
- Burst refresh: All rows refreshed immediately after one another
- Distributed refresh: Each row refreshed at a different time, at regular intervals

# DRAM Refresh (II)



Each pulse represents a refresh cycle

Required time to complete refresh of all rows

- Distributed refresh eliminates long pause times
- How else we can reduce the effect of refresh on performance?
  - Can we reduce the number of refreshes?

# Downsides of DRAM Refresh

-- Energy consumption: Each refresh consumes energy

-- Performance degradation: DRAM rank/bank unavailable while refreshed

-- QoS/predictability impact: (Long) pause times during refresh

-- Refresh rate limits DRAM density scaling



Liu et al., "RAIDR: Retention-aware Intelligent DRAM Refresh," ISCA 2012.

# 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 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

# 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 improve performance**
  - Reordering, row-buffer, bank, rank, bus management

- **Manage power consumption and thermals in DRAM**
  - Turn on/off DRAM chips, manage power modes

# DRAM Controller: Where to Place

- In chipset
  - + More flexibility to plug different DRAM types into the system
  - + Less power density in the CPU chip

- On CPU chip
  - + Reduced latency for main memory access
  - + Higher bandwidth between cores and controller
    - More information can be communicated (e.g. request's importance in the processing core)

# 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, older commands prioritized over younger ones

# DRAM Scheduling Policies (II)

- A scheduling policy is essentially a 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?

# 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

# 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 |

# 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 optimize for performance (in the presence of constraints)
  - Reordering is not simple
  - Predicting the future?

# 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.
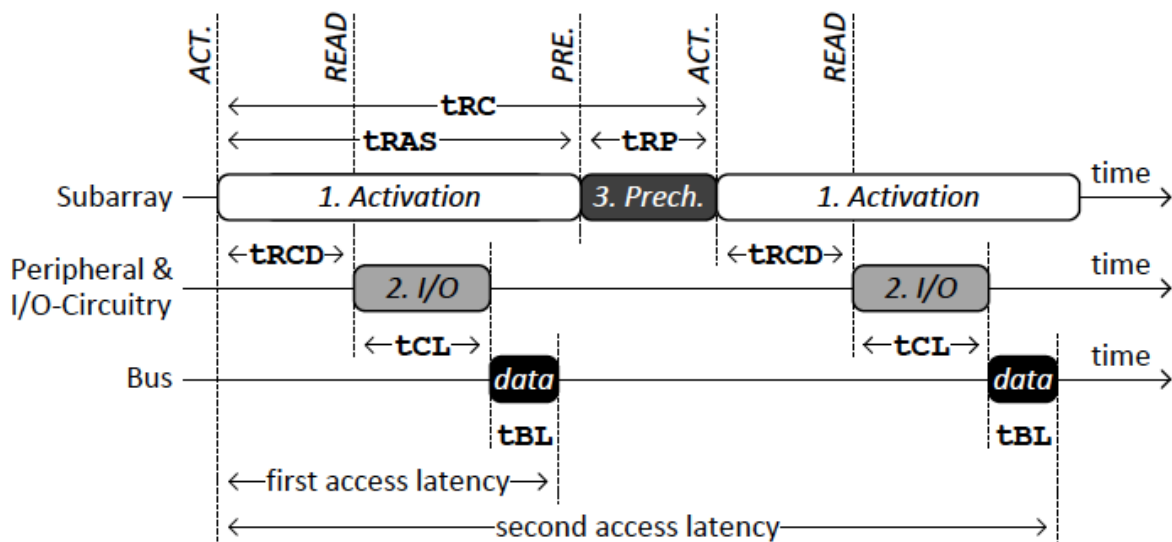


Figure 5. Three Phases of DRAM Access

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 |

We did not cover the remaining slides.

# Self-Optimizing DRAM Controllers

- Problem: DRAM controllers 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: Design a memory controller that adapts its scheduling policy decisions 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 that dynamically and continuously learns and employs the best scheduling policy.

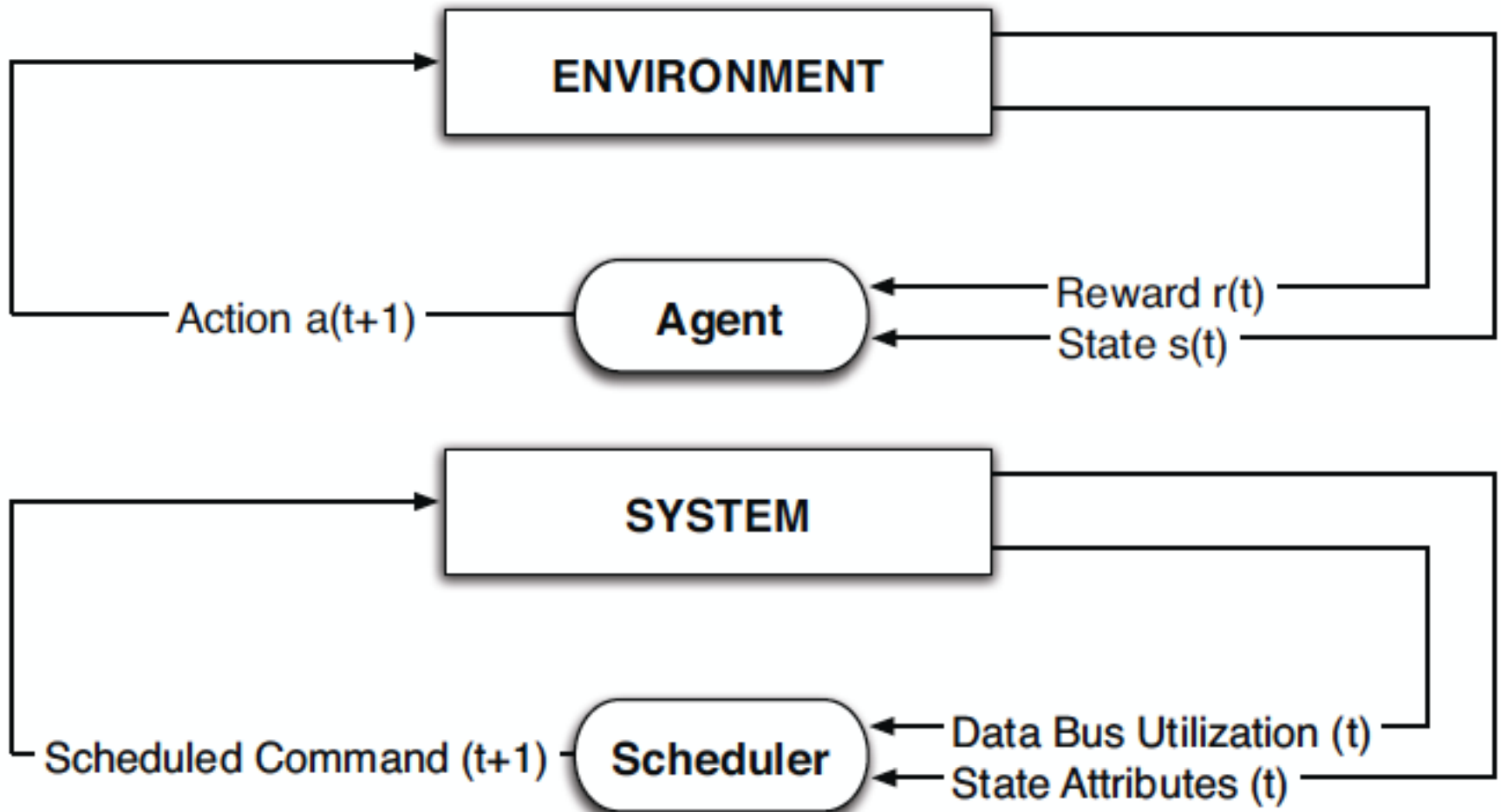# Self-Optimizing DRAM Controllers



**Figure 2:** (a) Intelligent agent based on reinforcement learning principles; (b) DRAM scheduler as an RL-agent

# Self-Optimizing DRAM Controllers

- Engin Ipek, <u>Onur Mutlu</u>, 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.
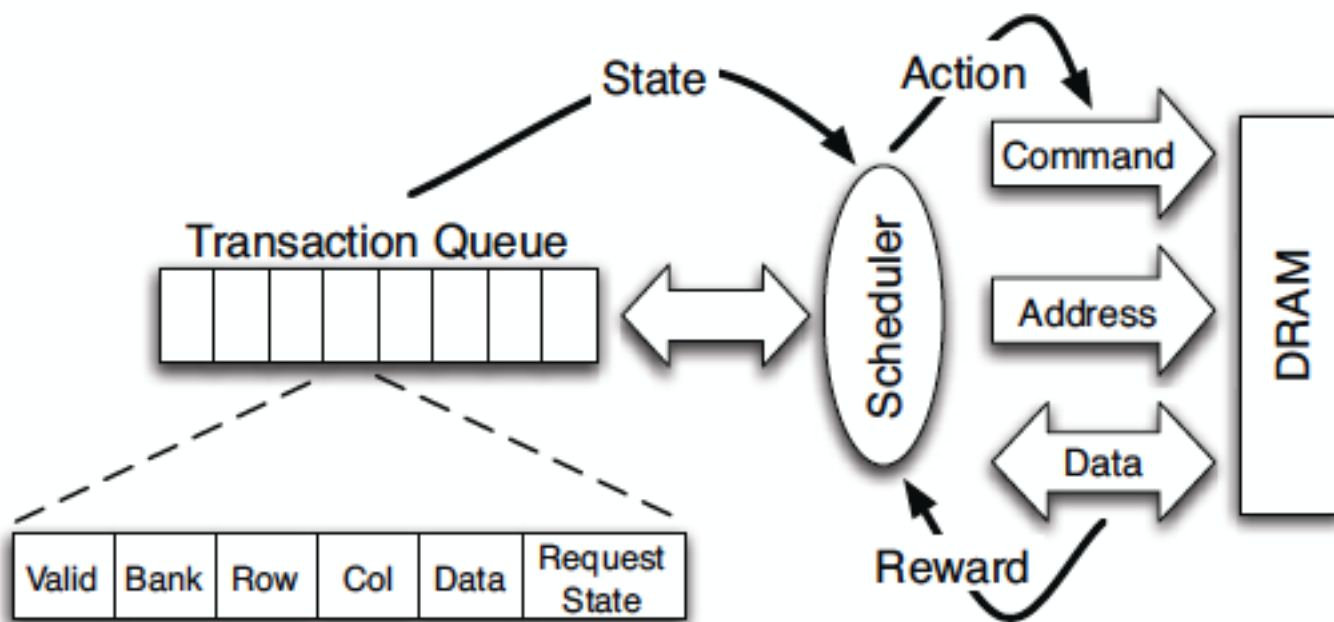


Figure 4: High-level overview of an RL-based scheduler.
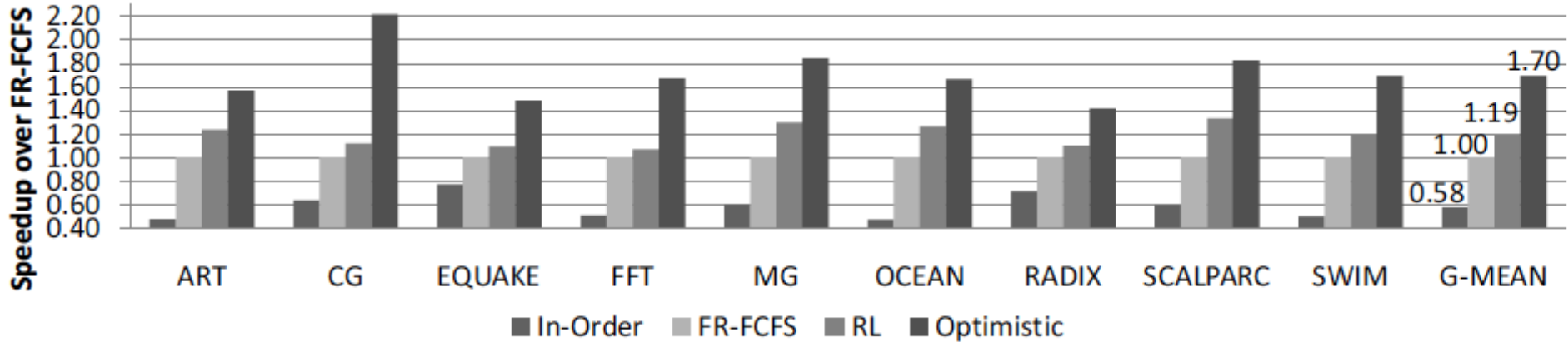
# Performance Results



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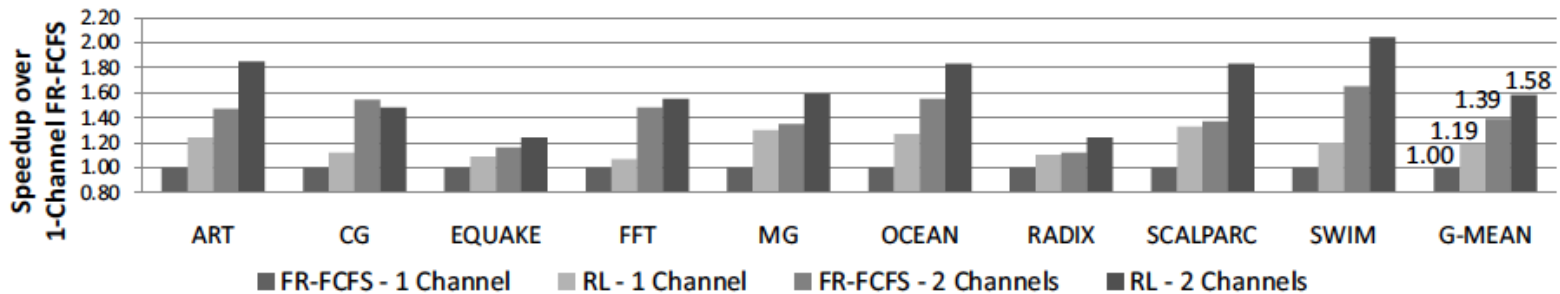


Figure 15: Performance comparison of FR-FCFS and RL-based memory controllers on systems with 6.4GB/s and 12.8GB/s peak DRAM bandwidth

# 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