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

**Carnegie Mellon**

*SAFARI*

# The Main Memory System



Processor and caches — Main Memory — Storage (SSD/HDD)
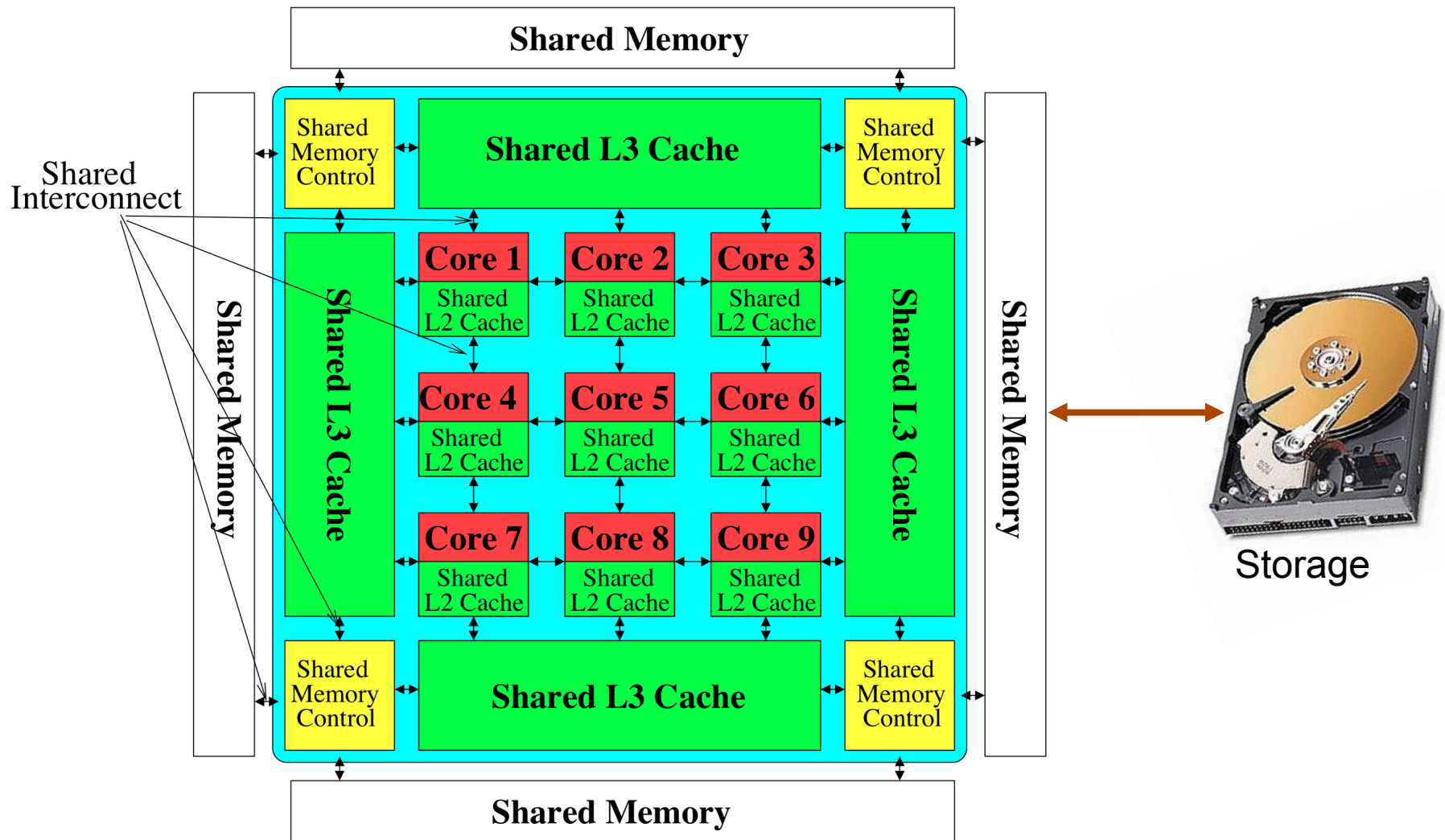
- **Main memory is a critical component of all computing systems**: server, mobile, embedded, desktop, sensor

- **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
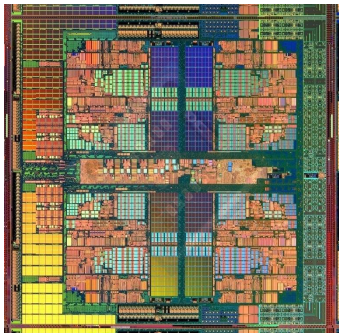
- DRAM technology scaling is ending
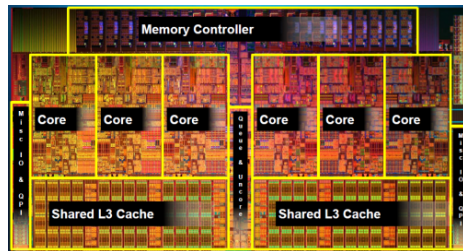
# Major Trends Affecting Main Memory (II)

- **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
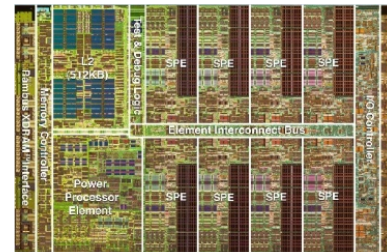
# Example Trend: Many Cores on Chip

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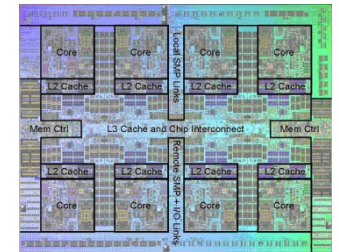- Large scale parallelism on chip
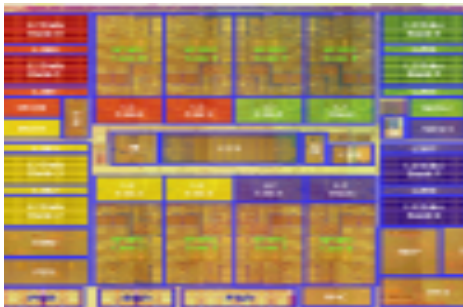


AMD Barcelona
4 cores

Intel Core i7
8 cores

IBM Cell BE
8+1 cores

IBM POWER7
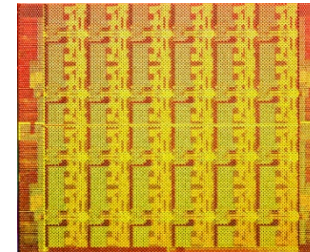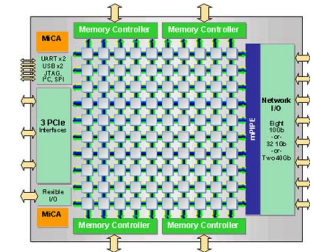8 cores

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

# 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.
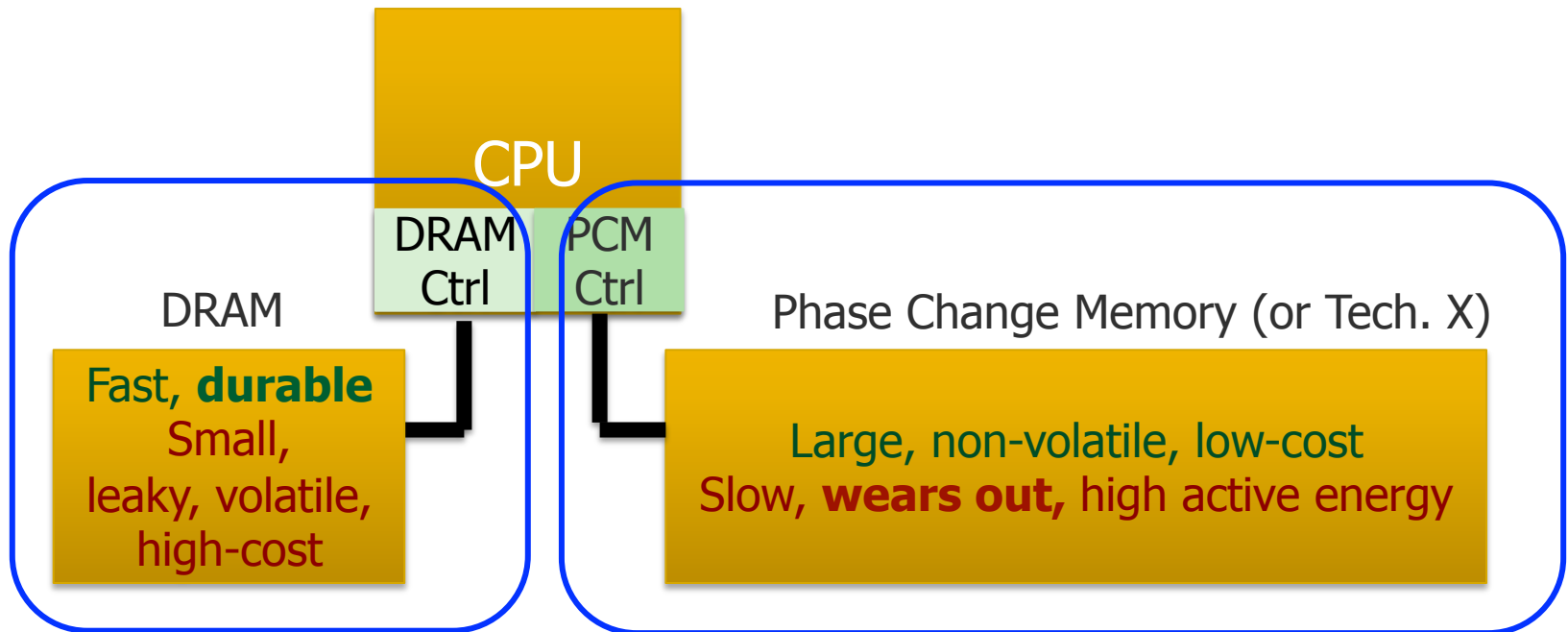
# Solution 2: Emerging Memory Technologies

- **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.

# Hybrid Memory Systems



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

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

# 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 Today

- **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

- Required 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)

# Readings and Videos

# Required 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.

# 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, Onur Mutlu, 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, Onur Mutlu, 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, Onur Mutlu, 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, Onur Mutlu, 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 Today

- 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 in the System

# Ideal Memory

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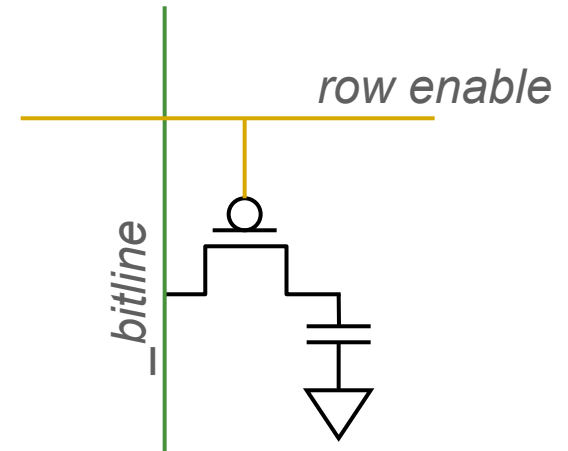- Infinite bandwidth (to support multiple accesses in parallel)

# The Problem

- Ideal memory's requirements oppose each other

- Bigger is slower
  - Bigger → 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

- 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.

_row enable_

_bitline_

# 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*

# 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

  - **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)

# 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.
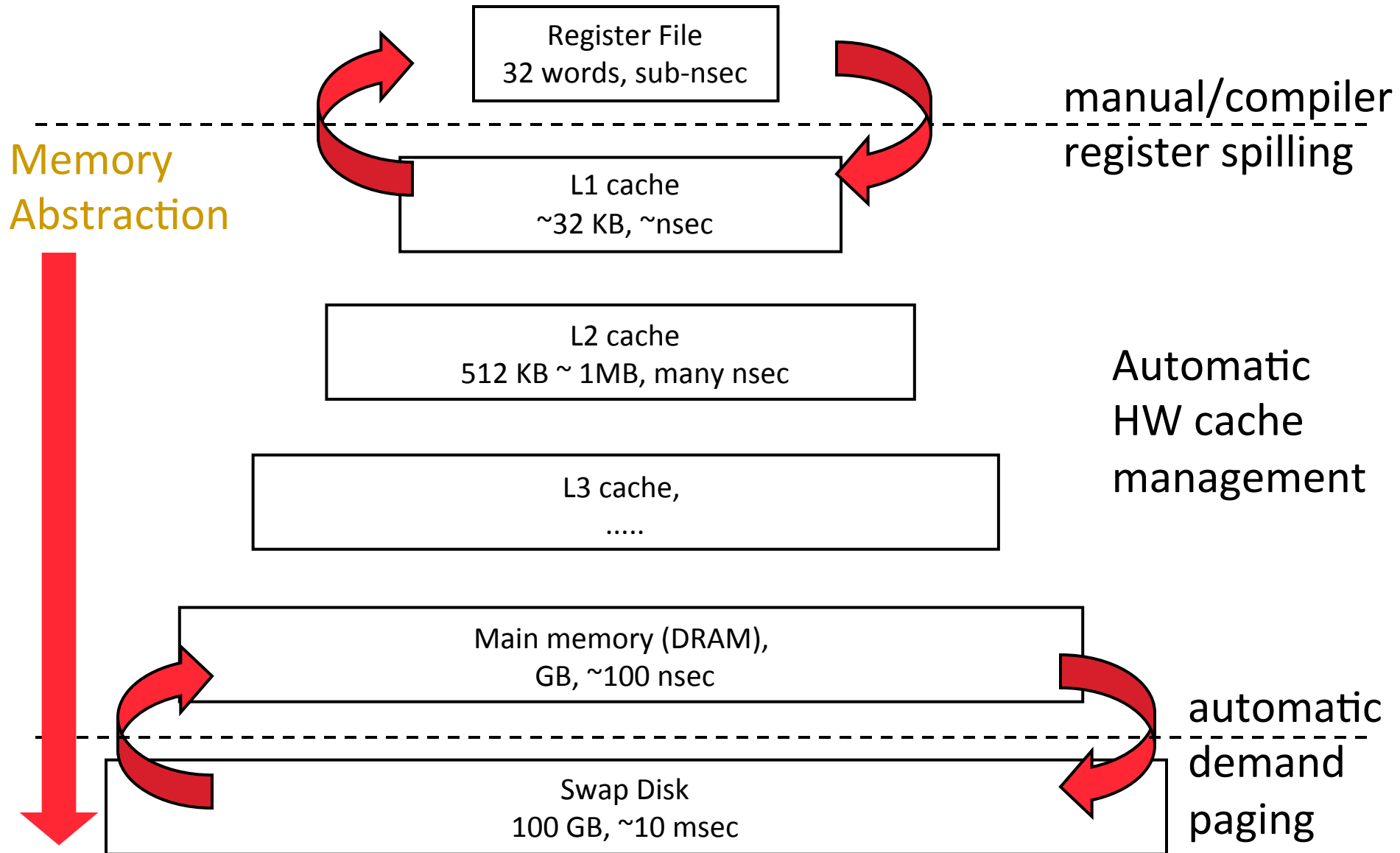    - "The use is discussed of a fast core memory of, say 32000 words as a slave to a slower core memory of, say, one million words in such a way that in practical cases the effective access time is nearer that of the fast memory than that of the slow memory."

# 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

Register File
32 words, sub-nsec

manual/compiler
register spilling

Memory
Abstraction

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

# The DRAM Subsystem

# 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 Buffer    CONFLICT ! HIT

Row 1

Column address 0 1 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



**"Mesh Topology"**

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

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

# 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

# DRAM Subsystem Organization

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

# The DRAM subsystem

# Breaking down a DIMM

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



Side view →
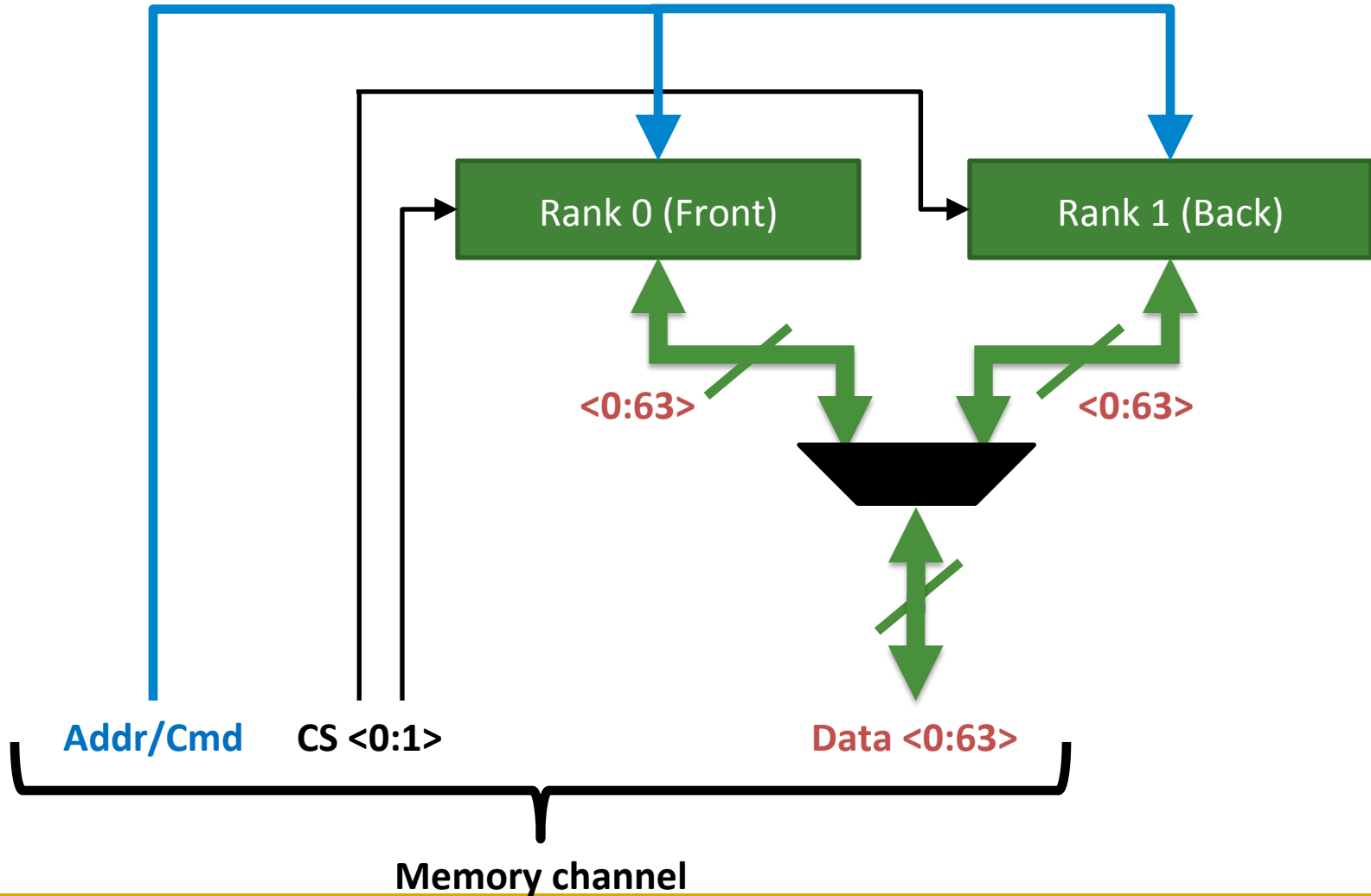
**SIDE**

4.00

**Front of DIMM**

SPD

**Back of DIMM**

# Breaking down a DIMM

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



Side view →

**SIDE**

4.00

**Front of DIMM**

**Back of DIMM**

SPD

**Rank 0:** collection of 8 chips

**Rank 1**

# Rank



Rank 0 (Front)   Rank 1 (Back)

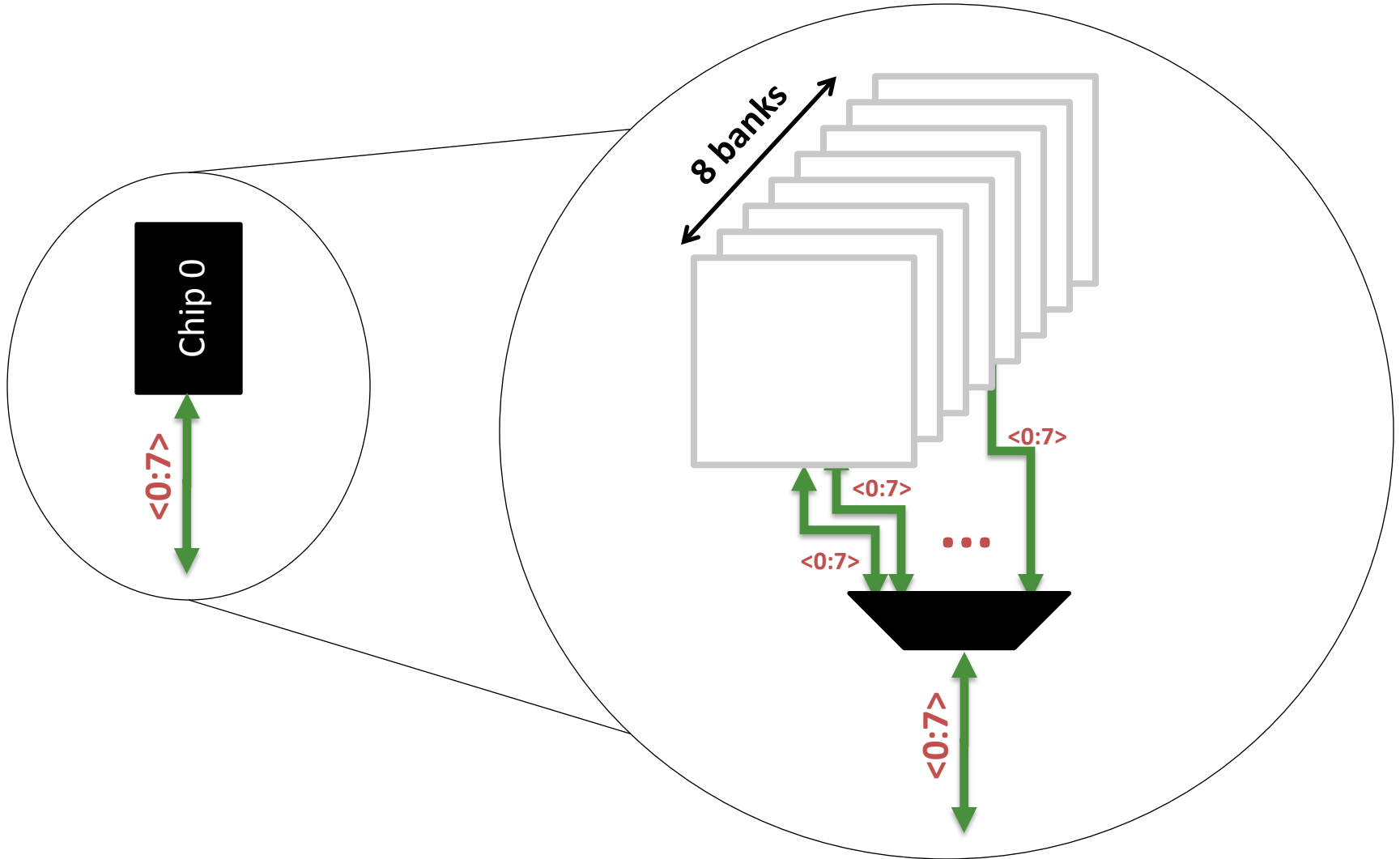<0:63>   <0:63>

Addr/Cmd   CS <0:1>   Data <0:63>

Memory channel

# Breaking down a Rank

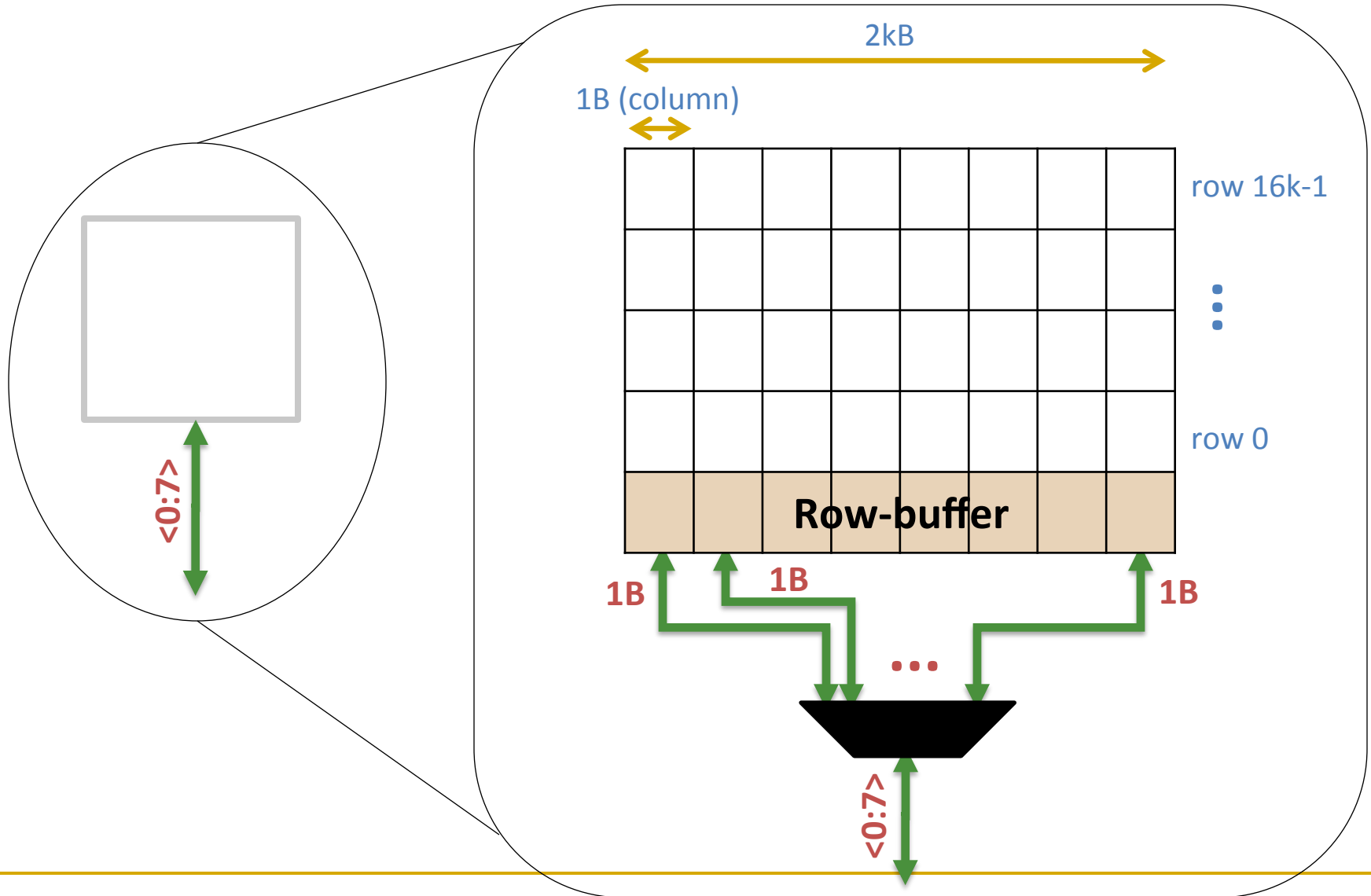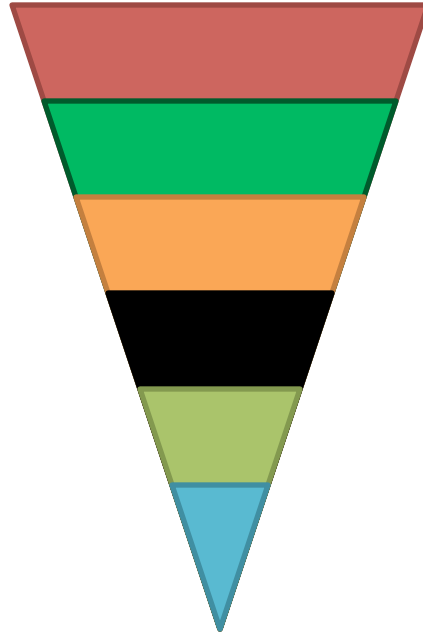# Breaking down a Chip

# 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
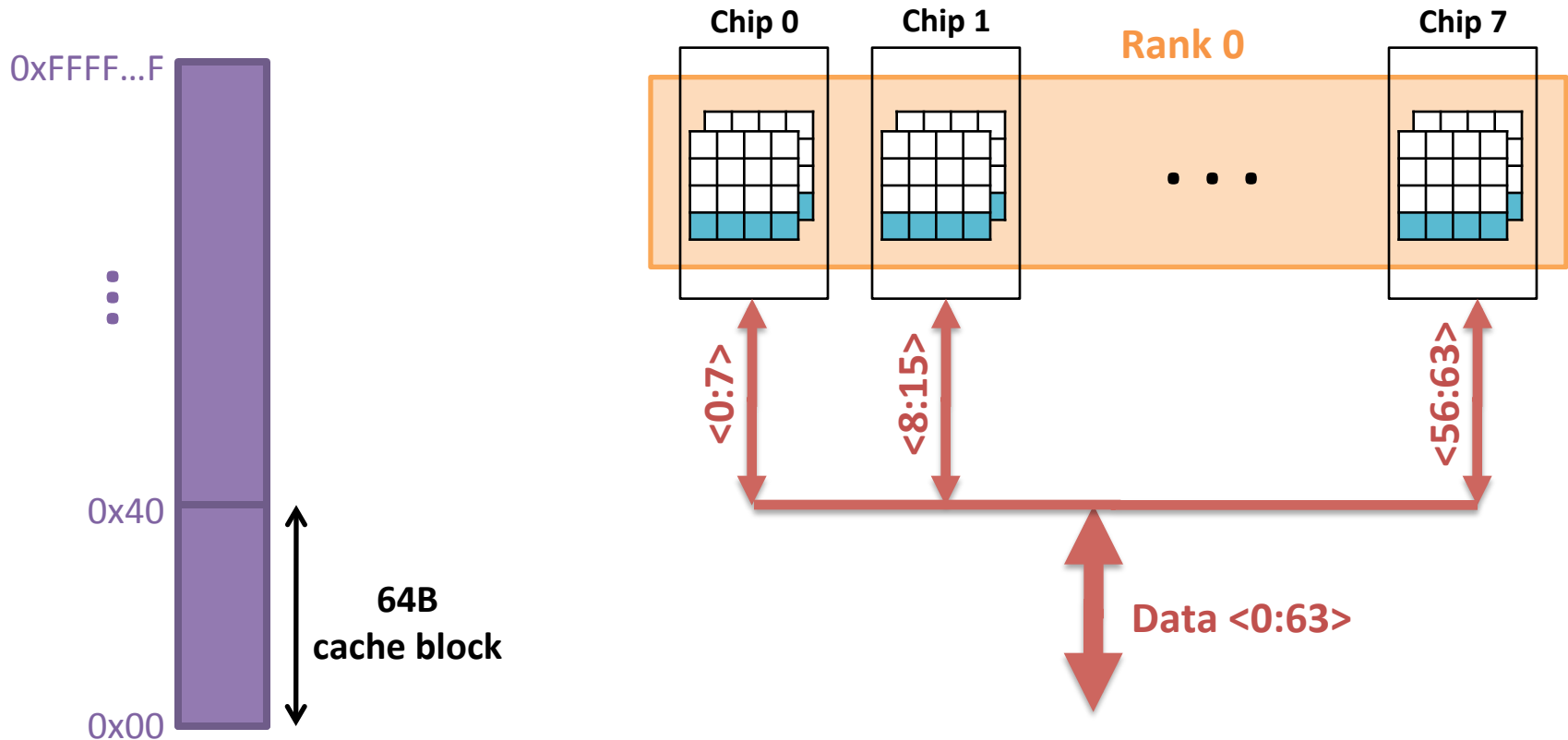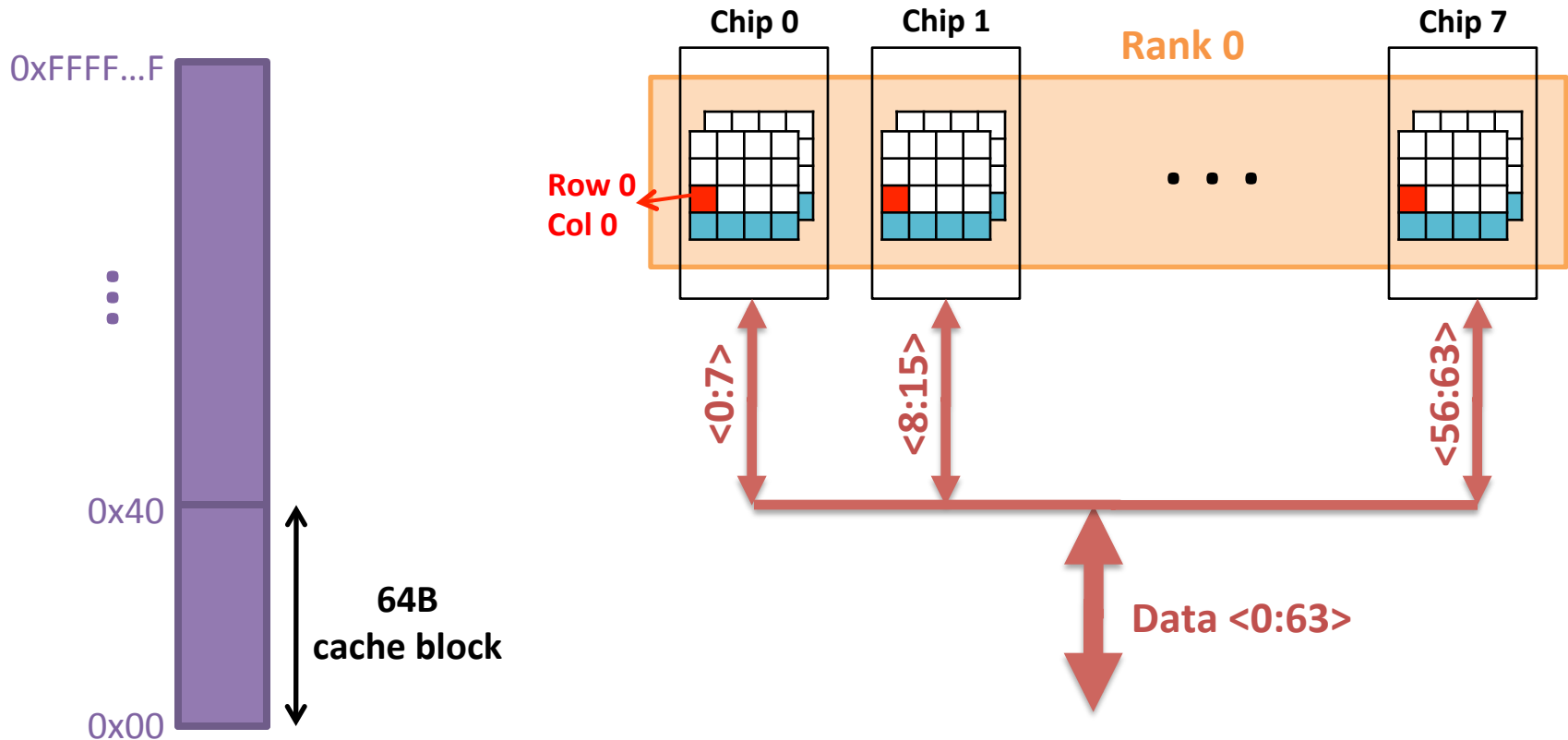
Channel 0

DIMM 0

Rank 0

intel Core™ i7

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block

**Physical memory space**

# Example: Transferring a cache block

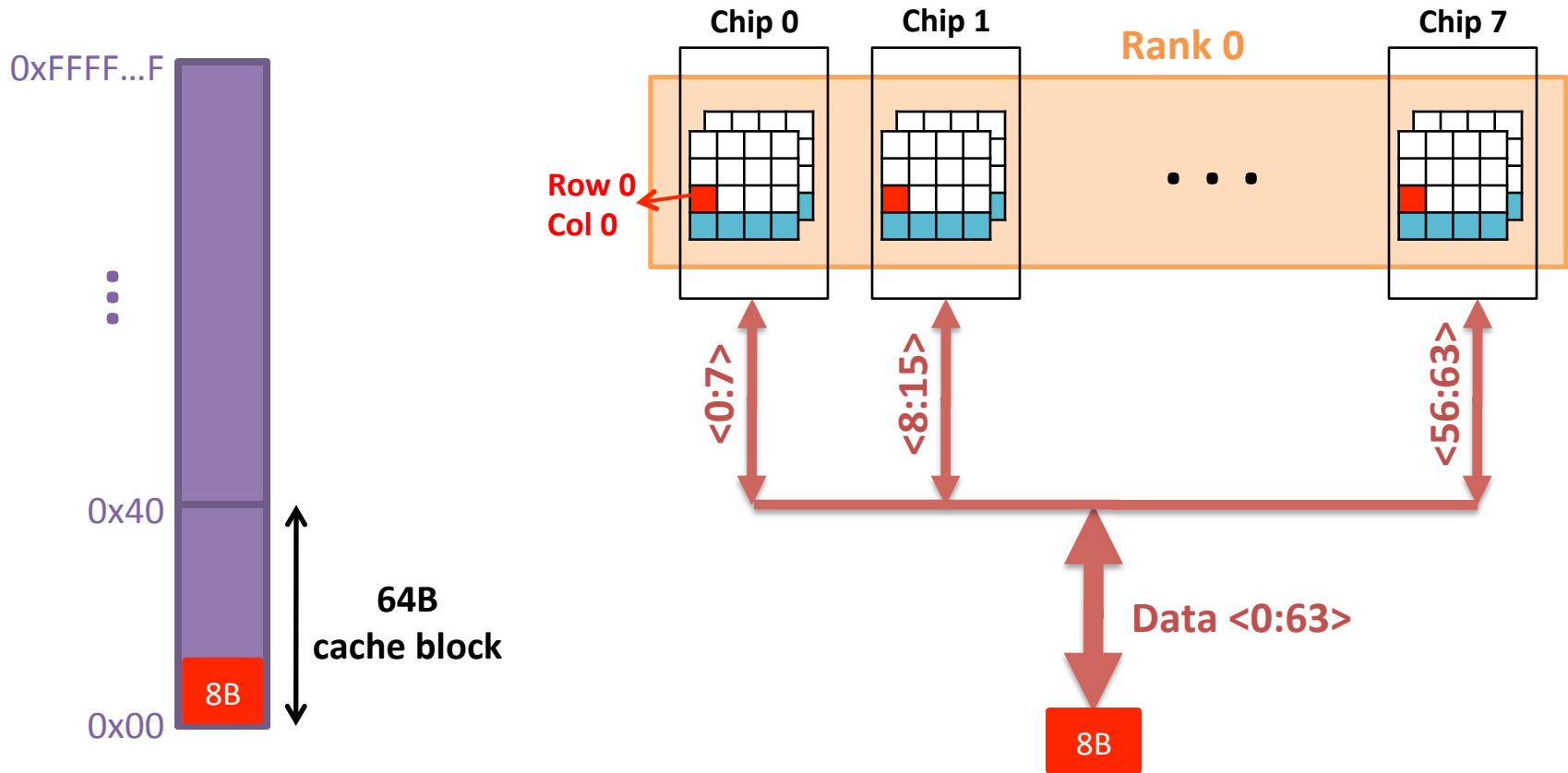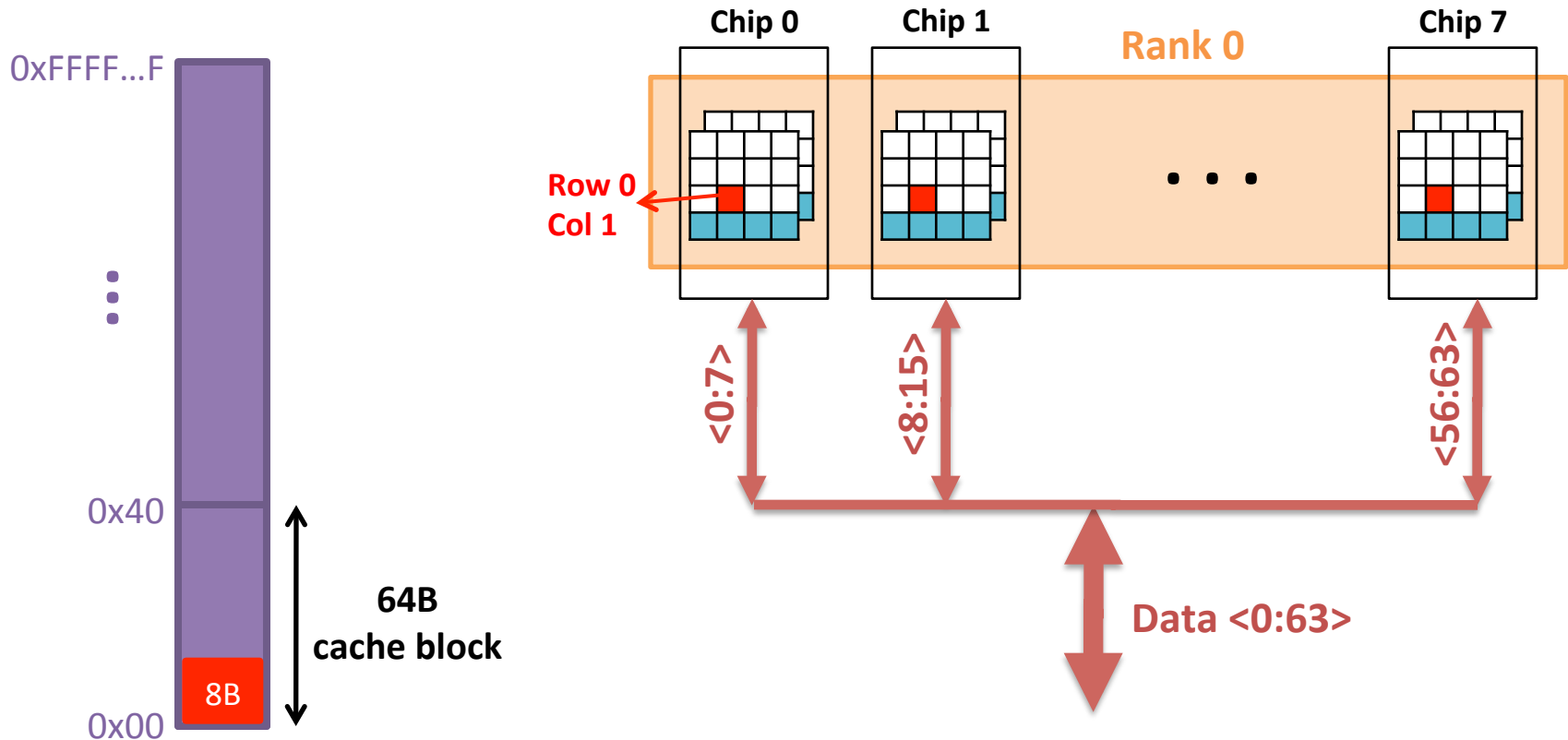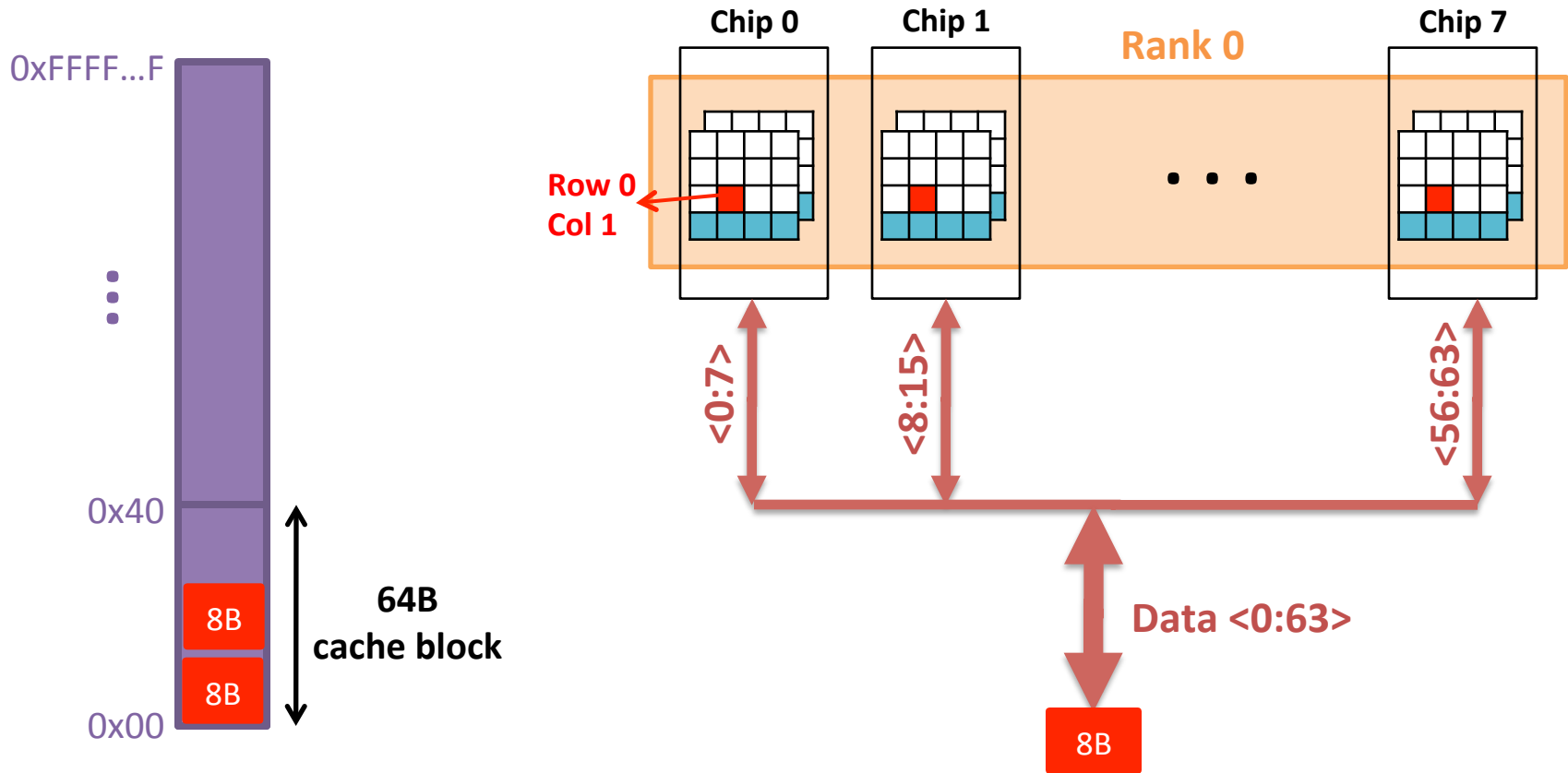**Physical memory space**

# 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 if row is "open" OR
  - RAS + CAS if array precharged OR
  - PRE + RAS + CAS (worst case)
- DRAM → CPU transfer time (through controller)

# Multiple Banks (Interleaving) and Channels

- **Multiple banks**
    - Enable <span style="color:red">concurrent DRAM accesses</span>
    - 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 <span style="color:red">separate data buses</span>
    - <span style="color:red">Increased bus bandwidth</span>

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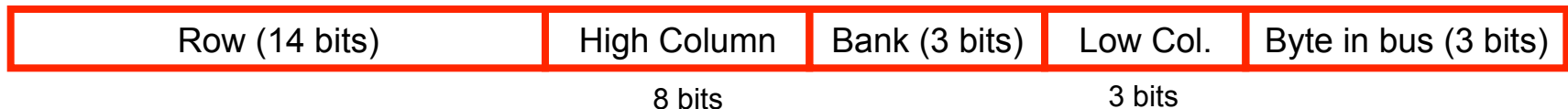    - Randomizing hash functions (XOR of different address bits)

# 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) |
|---|---|---|---|

- **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
- How about random accesses? Strided accesses?

# 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 |
|---|---|---|

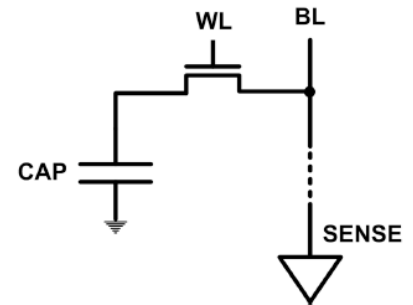| Row (14 bits) | Bank (3 bits) | Column (11 bits) | Byte in bus (3 bits) | PA |
|---|---|---|---|---|

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

- **It can perform page coloring to minimize bank conflicts**
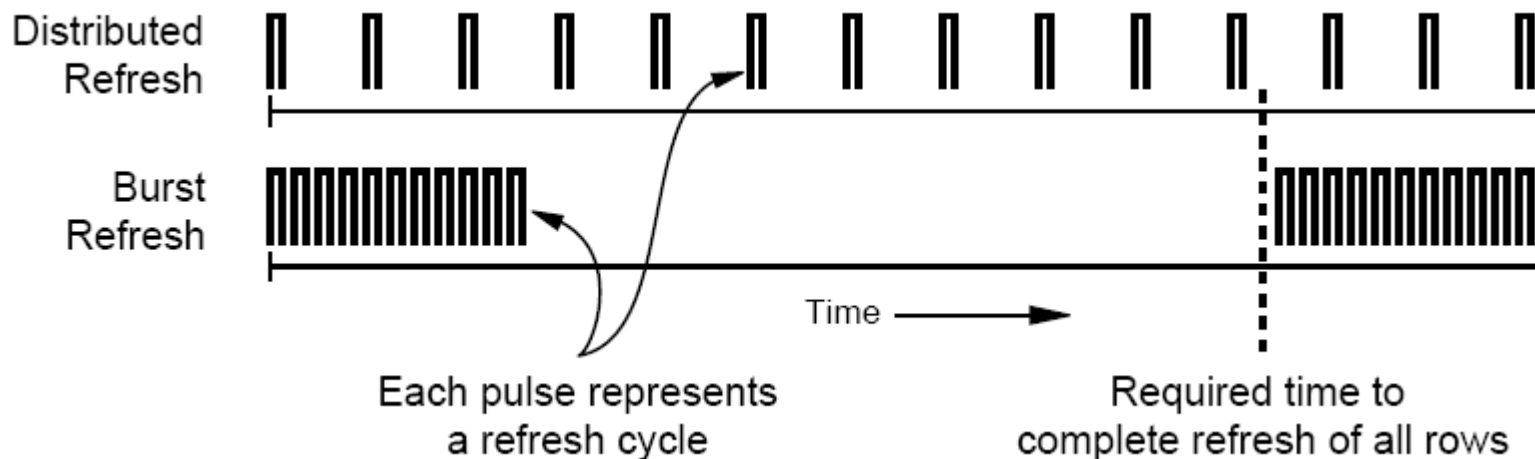- **Or to minimize inter-application interference**

# 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
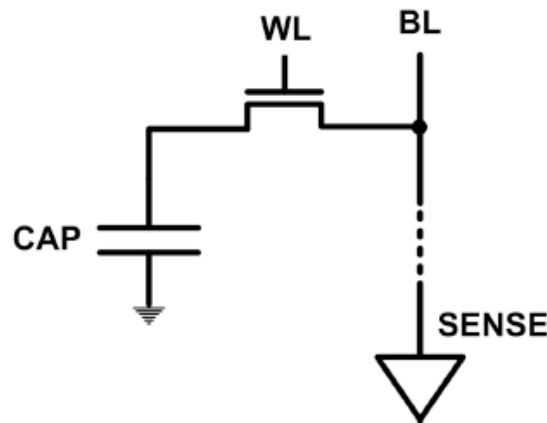
# DRAM Refresh (II)



- **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

- Downsides of 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

# 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
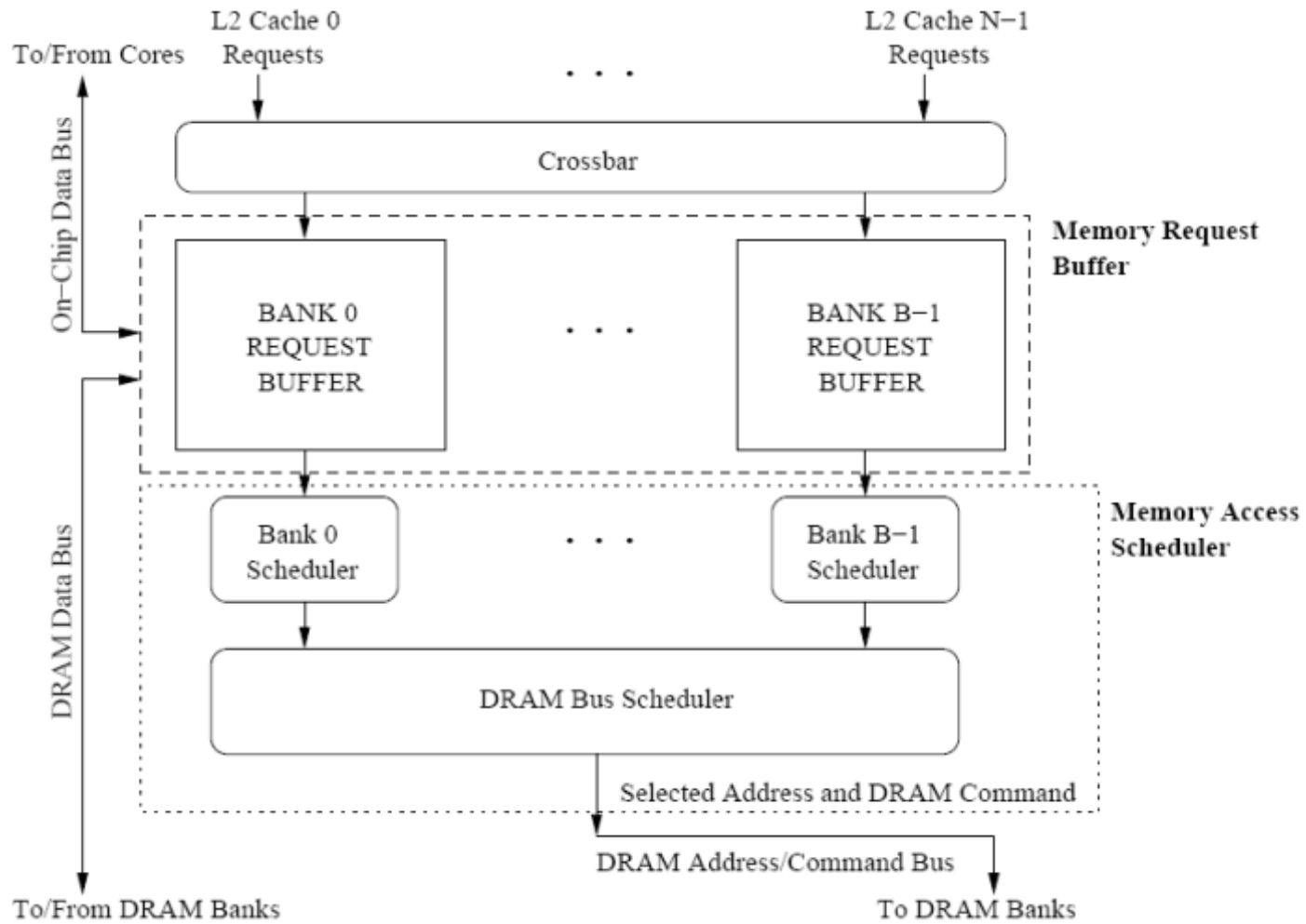
- **Ensure correct operation** of DRAM (refresh and timing)

- **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 |

87

# 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.
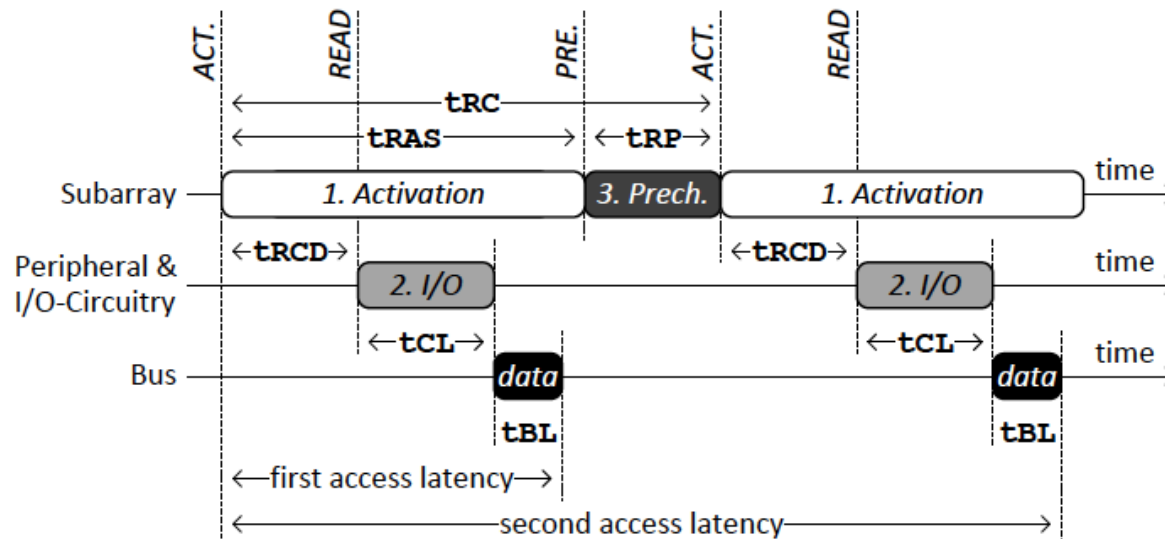


Figure 5. Three Phases of DRAM Access

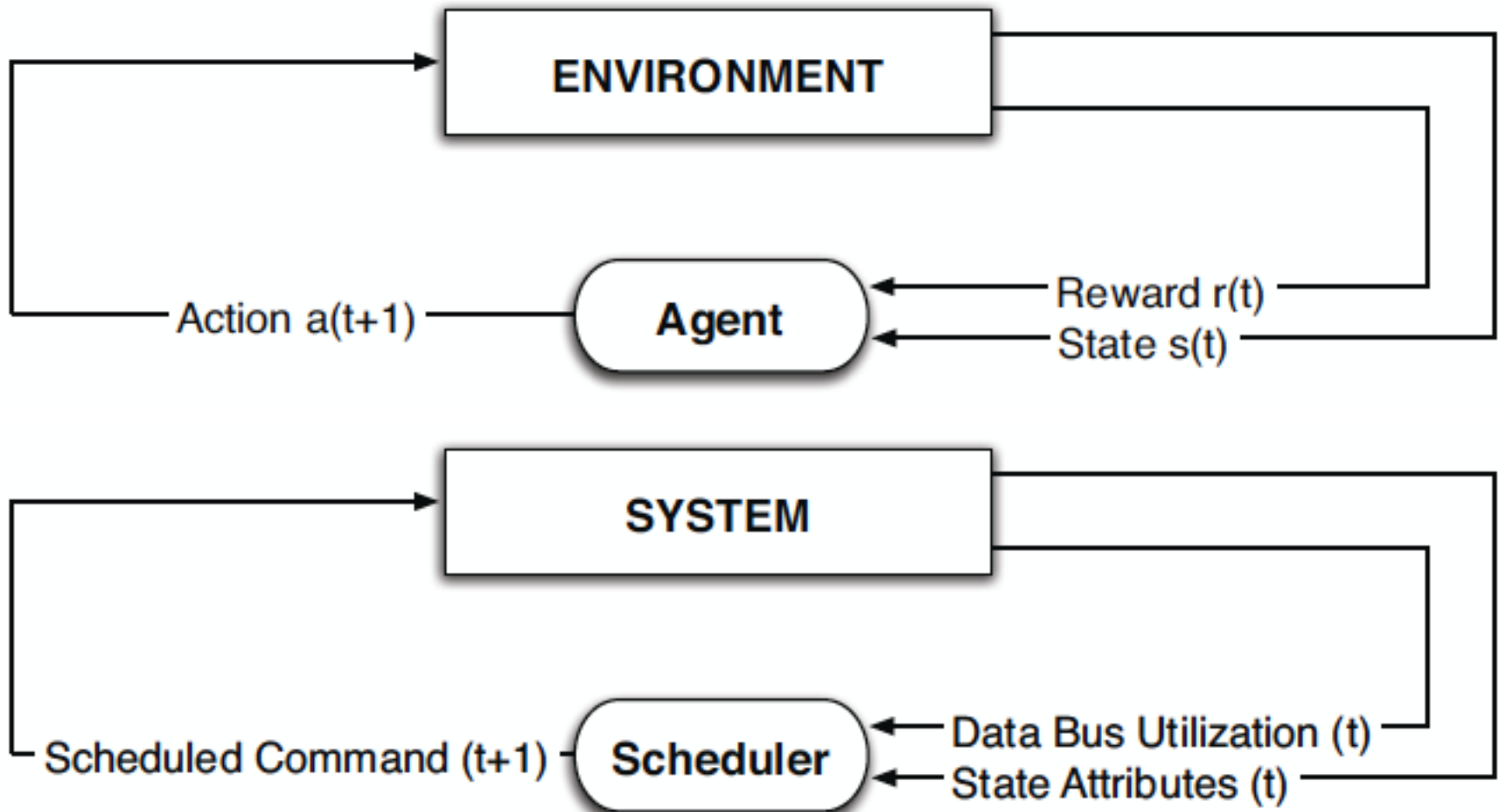Table 2. Timing Constraints (DDR3-1066) [43]

| 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 |

# 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, 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.

# 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)

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

# Agenda for Today

- 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

# 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.

# New DRAM Architectures

- **RAIDR: Reducing Refresh Impact**
- TL-DRAM: Reducing DRAM Latency
- SALP: Reducing Bank Conflict Impact
- RowClone: Fast Bulk Data Copy and Initialization

# RAIDR: Reducing DRAM Refresh Impact

Jamie Liu, Ben Jaiyen, Richard Veras, and Onur Mutlu,
**"RAIDR: Retention-Aware Intelligent DRAM Refresh"**
*Proceedings of the 39th International Symposium on Computer Architecture (**ISCA**),*
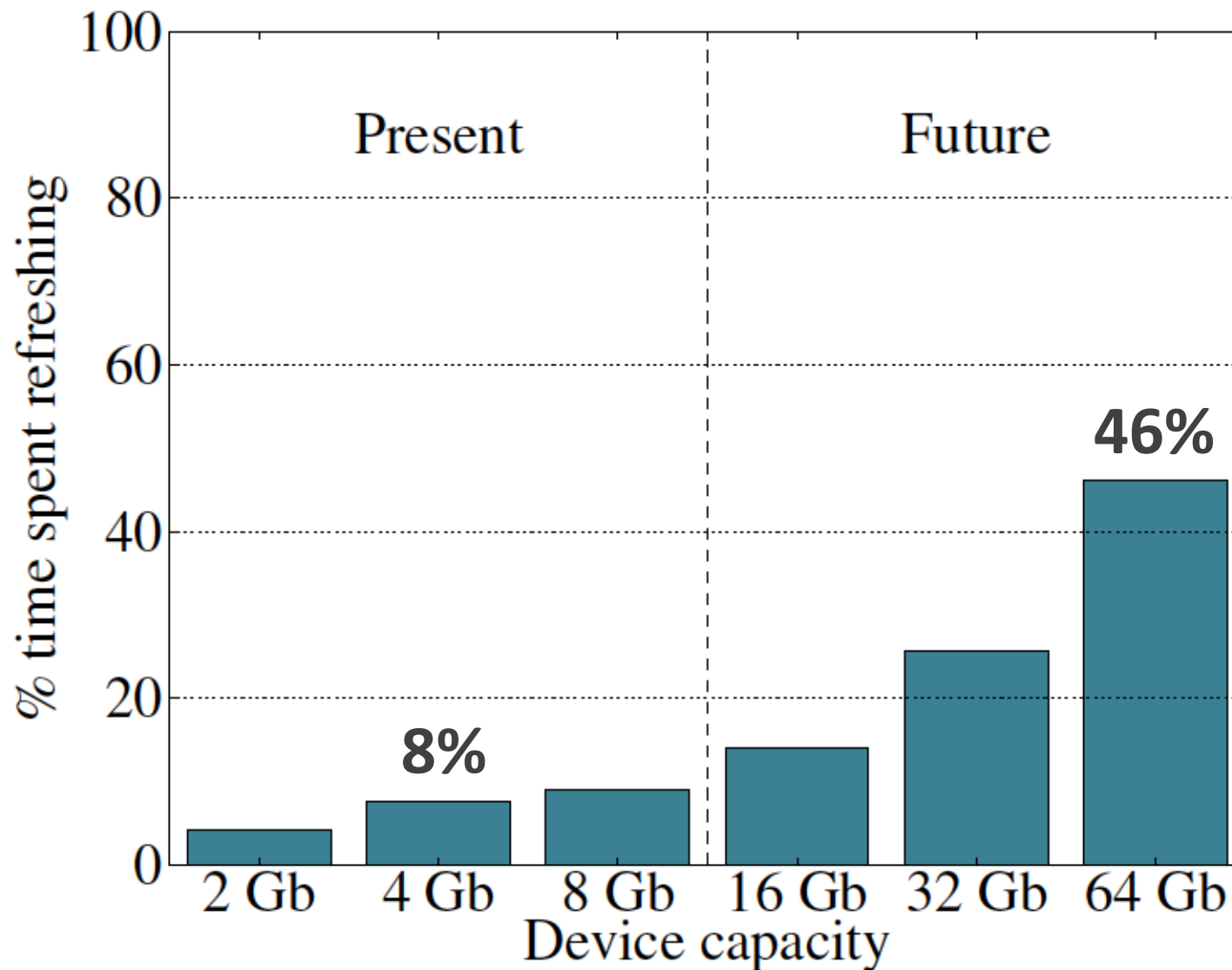Portland, OR, June 2012. Slides (pdf)

# DRAM Refresh

- DRAM capacitor charge leaks over time

- The memory controller needs to refresh each row periodically to restore charge
  - Activate + precharge each row every N ms
  - Typical N = 64 ms

- Downsides of 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

# Refresh Today: Auto Refresh



Columns

BANK 0   BANK 1   BANK 2   BANK 3

Rows

Row Buffer

DRAM Bus

DRAM CONTROLLER

A batch of rows are periodically refreshed via the auto-refresh command

# Refresh Overhead: Performance

# Refresh Overhead: Energy

SAFARI

# Problem with Conventional Refresh

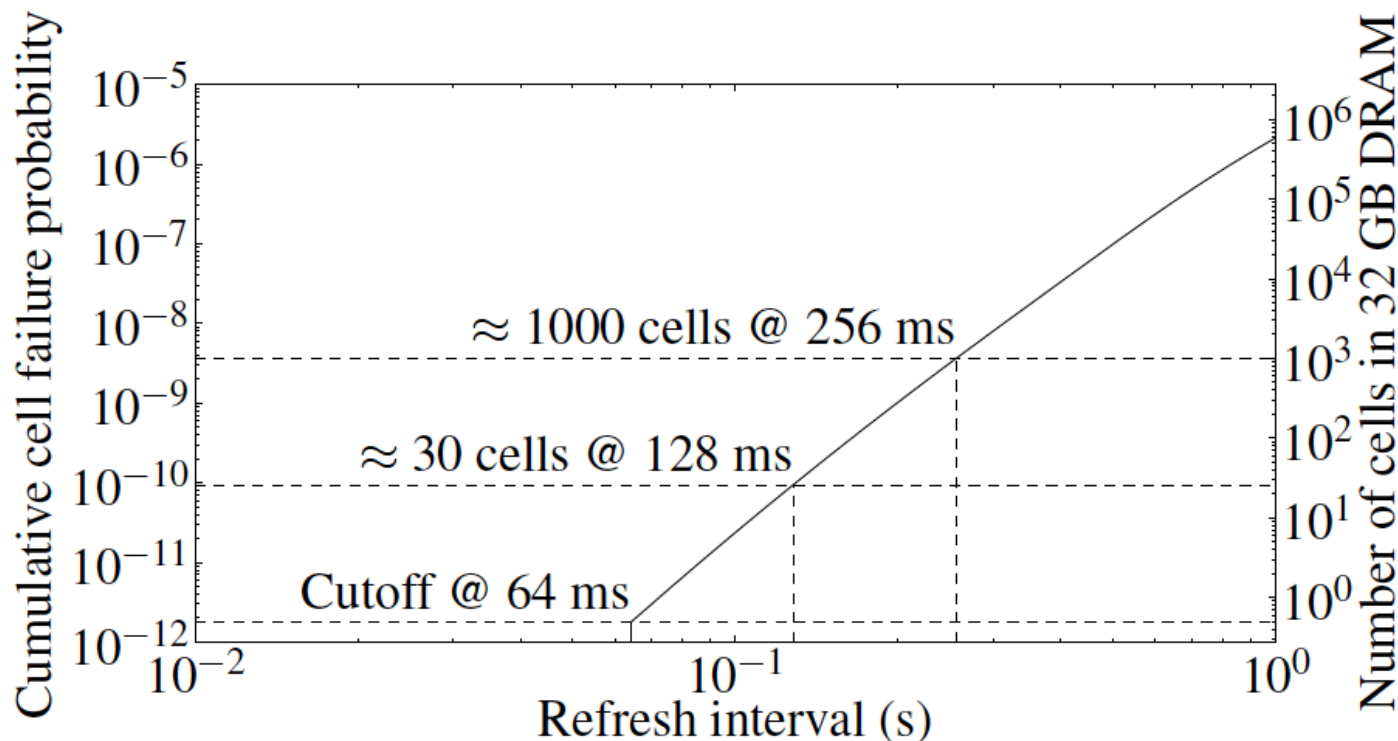- Today: Every row is refreshed at the same rate



- Observation: Most rows can be refreshed much less often without losing data [Kim+, EDL'09]
- Problem: No support in DRAM for different refresh rates per row

**SAFARI**

# Retention Time of DRAM Rows

- Observation: Only very few rows need to be refreshed at the worst-case rate



- Can we exploit this to reduce refresh operations at low cost?

# Reducing DRAM Refresh Operations

- **Idea:** Identify the retention time of different rows and refresh each row at the frequency it needs to be refreshed

- **(Cost-conscious) Idea:** Bin the rows according to their minimum retention times and refresh rows in each bin at the refresh rate specified for the bin
  - e.g., a bin for 64-128ms, another for 128-256ms, ...

- **Observation:** Only very few rows need to be refreshed very frequently [64-128ms] → Have only a few bins → Low HW overhead to achieve large reductions in refresh operations

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

# RAIDR: Mechanism

64-128ms

>256ms

1.25KB storage in controller for 32GB DRAM memory

128-256ms

bins at different rates

→ probe Bloom Filters to determine refresh rate of a row

# 1. Profiling

To profile a row:

1. Write data to the row
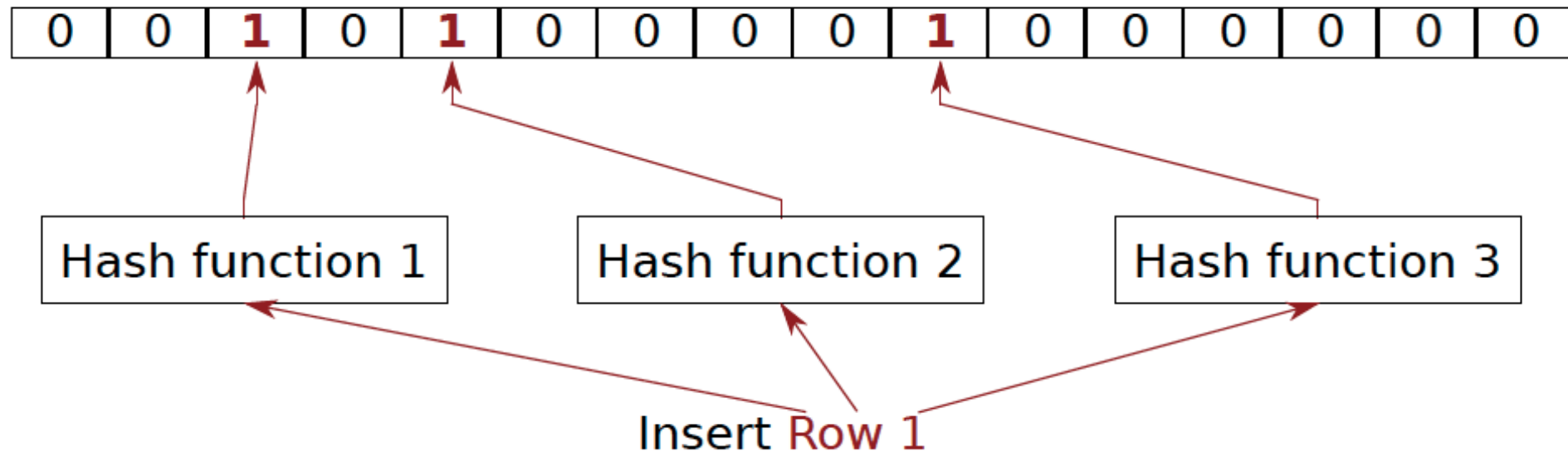2. Prevent it from being refreshed
3. Measure time before data corruption

|  | Row 1 | Row 2 | Row 3 |
|---|---|---|---|
| Initially | 11111111… | 11111111… | 11111111… |
| After 64 ms | 11111111… | 11111111… | 11111111… |
| After 128 ms | 11011111… (64–128ms) | 11111111… | 11111111… |
| After 256 ms |  | 11111011… (128–256ms) | 11111111… (>256ms) |

# 2. Binning

- **How to efficiently and scalably store rows into retention time bins?**

- Use Hardware Bloom Filters [Bloom, CACM 1970]

Example with 64–128ms bin:

| 0 | 0 | **1** | 0 | **1** | 0 | 0 | 0 | 0 | **1** | 0 | 0 | 0 | 0 | 0 | 0 |

Hash function 1    Hash function 2    Hash function 3

Insert Row 1

# Benefits of Bloom Filters as Bins

- **False positives:** a row may be declared present in the Bloom filter even if it was never inserted
  - **Not a problem:** Refresh some rows more frequently than needed

- **No false negatives:** rows are never refreshed less frequently than needed (no correctness problems)

- **Scalable:** a Bloom filter never overflows (unlike a fixed-size table)

- **Efficient:** No need to store info on a per-row basis; simple hardware → 1.25 KB for 2 filters for 32 GB DRAM system

**SAFARI**

# 3. Refreshing (RAIDR Refresh Controller)

Choose a refresh candidate row

↓

Determine which bin the row is in

↓

Determine if refreshing is needed

# 3. Refreshing (RAIDR Refresh Controller)

Memory controller
chooses each row
as a refresh candidate
every 64ms

↓

Row in 64-128ms bin? → Row in 128-256ms bin?
(First Bloom filter: 256B)    (Second Bloom filter: 1KB)

↓                    ↓                    ↓

Refresh the row     Every other 64ms window,     Every 4th 64ms window,
                    refresh the row              refresh the row

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

# Tolerating Temperature Changes

- Change in temperature causes retention time of all cells to change by a uniform and predictable factor

- Refresh rate scaling: increase the refresh rate for all rows uniformly, depending on the temperature

- Implementation: counter with programmable period
  - Lower temperature $\Rightarrow$ longer period $\Rightarrow$ less frequent refreshes
  - Higher temperature $\Rightarrow$ shorter period $\Rightarrow$ more frequent refreshes

# RAIDR Results

- Baseline:
  - 32 GB DDR3 DRAM system (8 cores, 512KB cache/core)
  - 64ms refresh interval for all rows

- RAIDR:
  - 64–128ms retention range: 256 B Bloom filter, 10 hash functions
  - 128–256ms retention range: 1 KB Bloom filter, 6 hash functions
  - Default refresh interval: 256 ms

- Results on SPEC CPU2006, TPC-C, TPC-H benchmarks
  - 74.6% refresh reduction
  - ~16%/20% DRAM dynamic/idle power reduction
  - ~9% performance improvement

# RAIDR Refresh Reduction



32 GB DDR3 DRAM system

# RAIDR: Performance



RAIDR performance benefits increase with workload's memory intensity

# RAIDR: DRAM Energy Efficiency



RAIDR energy benefits increase with memory idleness

**SAFARI**

# DRAM Device Capacity Scaling: Performance



RAIDR performance benefits increase with DRAM chip capacity

# DRAM Device Capacity Scaling: Energy



RAIDR energy benefits increase with DRAM chip capacity

**SAFARI**

# New DRAM Architectures

- RAIDR: Reducing Refresh Impact

- TL-DRAM: Reducing DRAM Latency

- SALP: Reducing Bank Conflict Impact

- RowClone: Fast Bulk Data Copy and Initialization

# Tiered-Latency DRAM: Reducing DRAM Latency

Donghyuk Lee, Yoongu Kim, Vivek Seshadri, Jamie Liu, Lavanya Subramanian, and Onur Mutlu,
**"Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture"**

# Historical DRAM Latency-Capacity Trend



*DRAM latency continues to be a critical bottleneck*

# What Causes the Long Latency?

**DRAM Chip**



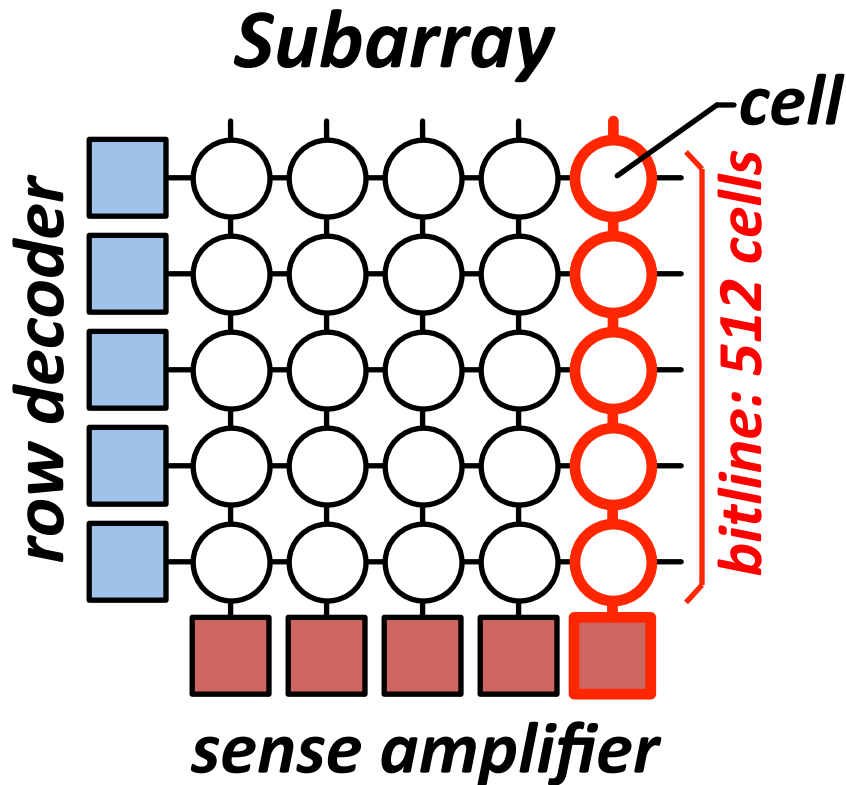*subarray*

*I/O*

*channel*

*subarray*

*cell*

wordline

*row decoder*

capacitor

access transistor

bitline

*sense amplifier*

123

# What Causes the Long Latency?

**DRAM Chip**

**subarray**

**row decoder**

**subarray**

**sense amplifier**

**row addr.**

**Subarray**

**I/O**

**I/O**

**channel**

**column addr.**

**mux**

DRAM Latency = ~~Subarray Latency~~ + ~~I/O Latency~~

**Dominant**

# Why is the Subarray So Slow?

**Subarray**

*row decoder*

cell

bitline: 512 cells

sense amplifier

**Cell**

*row decoder*

wordline

capacitor

access transistor

bitline

*sense amplifier*

**large sense amplifier**

- **Long bitline**
  - **Amortizes sense amplifier cost → Small area**
  - **Large bitline capacitance → High latency & power**

125

# Trade-Off: Area (Die Size) vs. Latency

**Long Bitline**

**Short Bitline**

**Faster**

**Smaller**

**Trade-Off: Area vs. Latency**

# Trade-Off: Area (Die Size) vs. Latency

# Approximating the Best of Both Worlds

**Long Bitline**

**Our Proposal**

**Short Bitline**

*Small Area*

~~*Large Area*~~

~~*High Latency*~~

*Low Latency*

*Need Isolation*

*Add Isolation Transistors*

*tline → Fast*

# Approximating the Best of Both Worlds

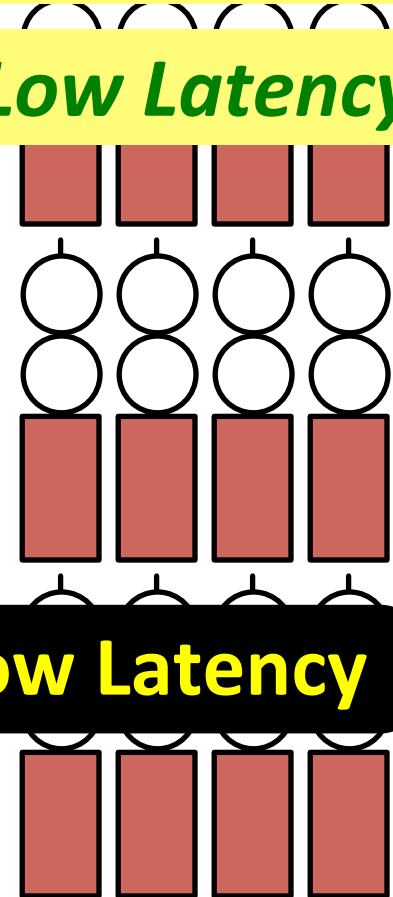**Long Bitline** **Tiered-Latency DRAM** **Short Bitline**

Small Area | Small Area | ~~Large Area~~
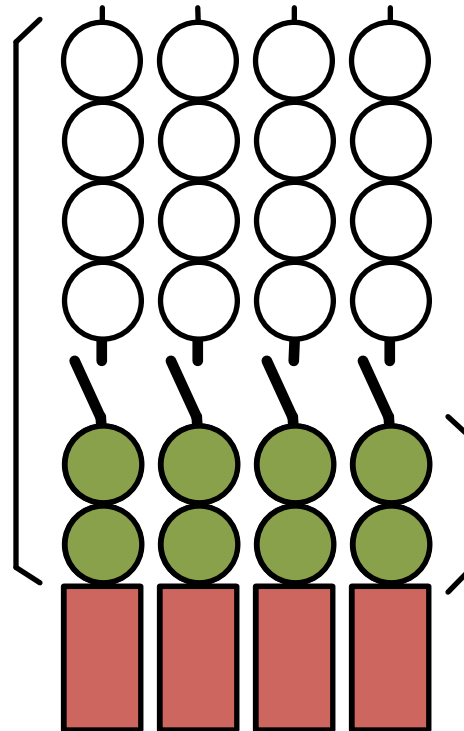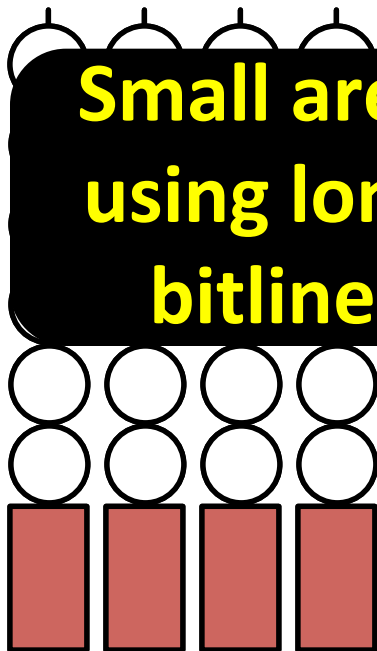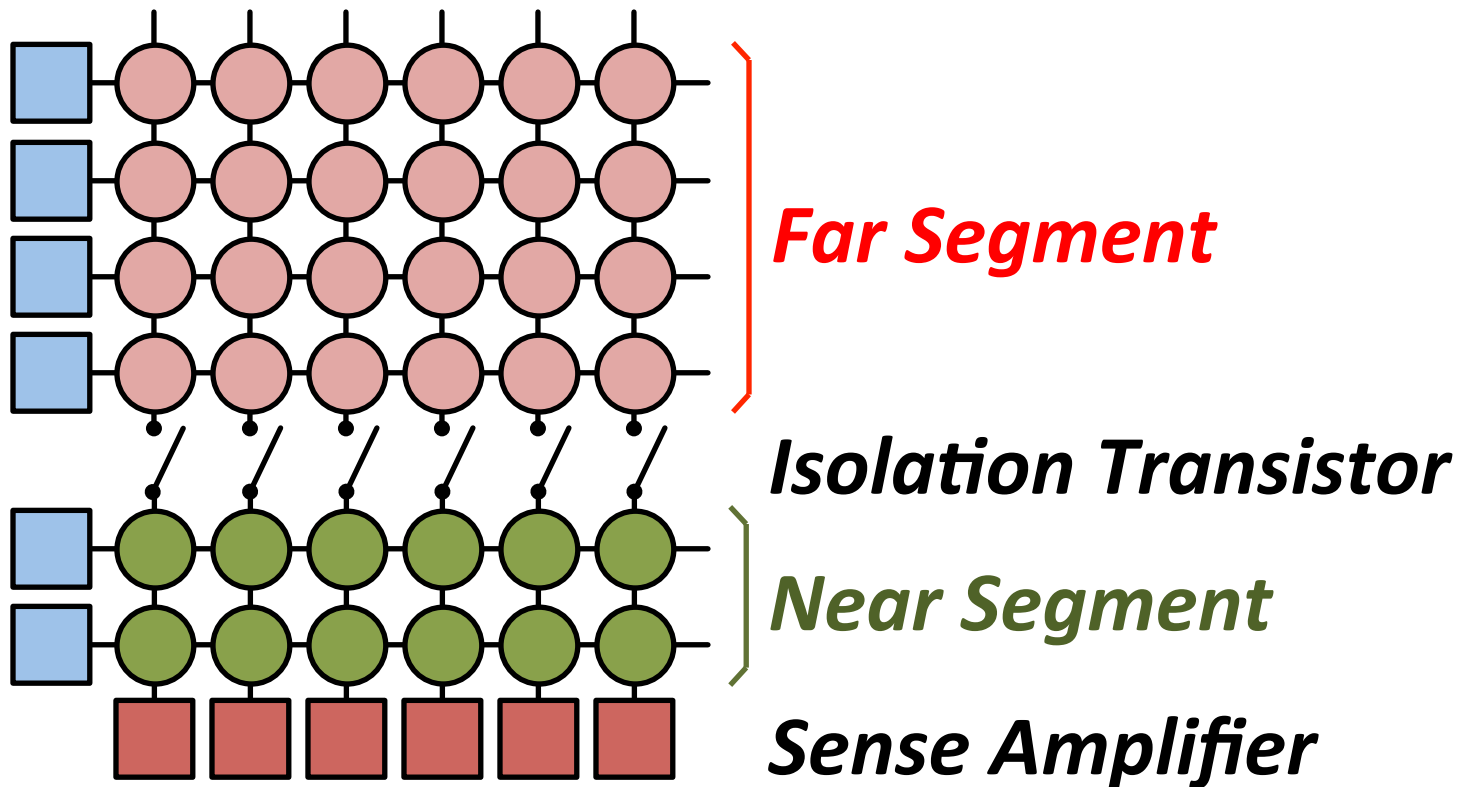
~~High Latency~~ | Low Latency | Low Latency

**Small area using long bitline**

**Low Latency**

# Tiered-Latency DRAM

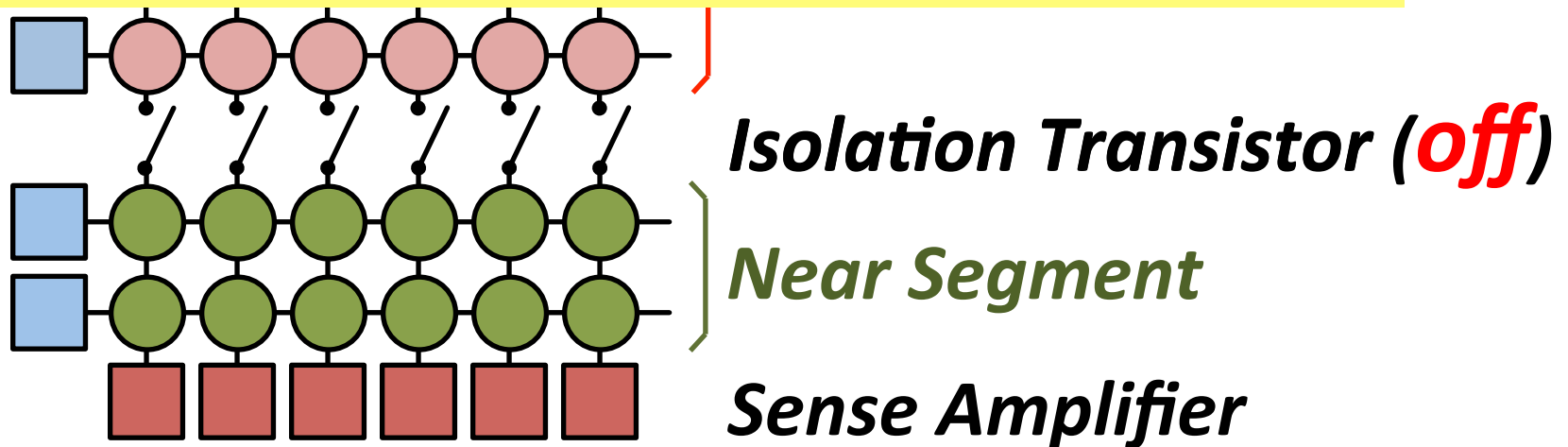- Divide a bitline into two segments with an **isolation transistor**

*Far Segment*

*Isolation Transistor*

*Near Segment*

*Sense Amplifier*

# Near Segment Access

- **Turn *off* the isolation transistor**

**Reduced bitline length**

**Reduced bitline capacitance**

→ **Low latency & low power**

*Isolation Transistor (off)*

*Near Segment*

*Sense Amplifier*
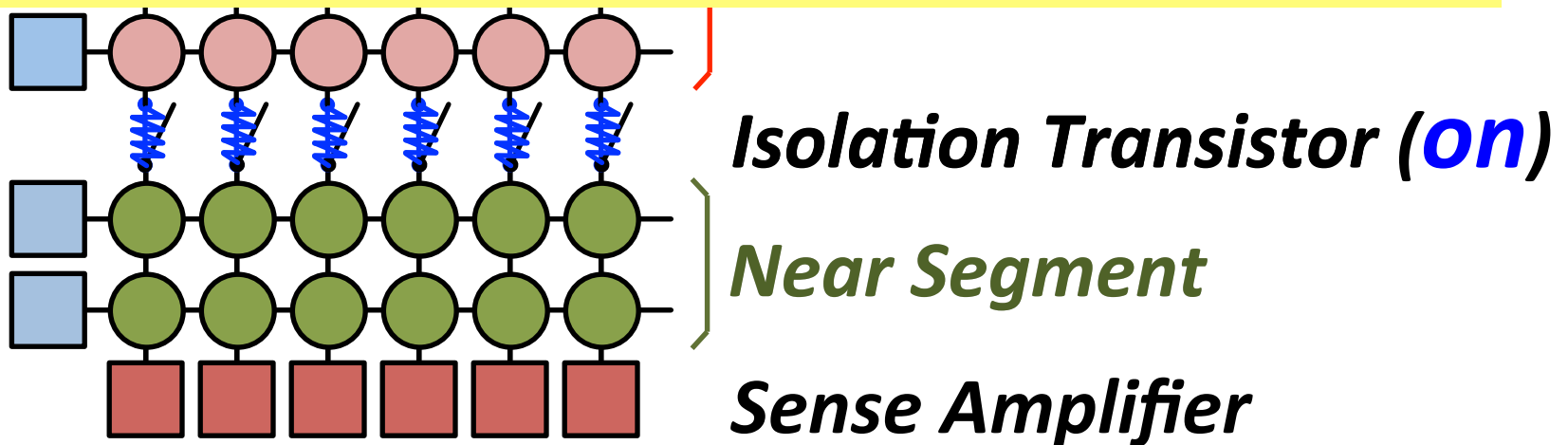
# Far Segment Access

- Turn *on* the isolation transistor

Long bitline length

Large bitline capacitance

Additional resistance of isolation transistor

➔ High latency & high power

*Isolation Transistor (on)*
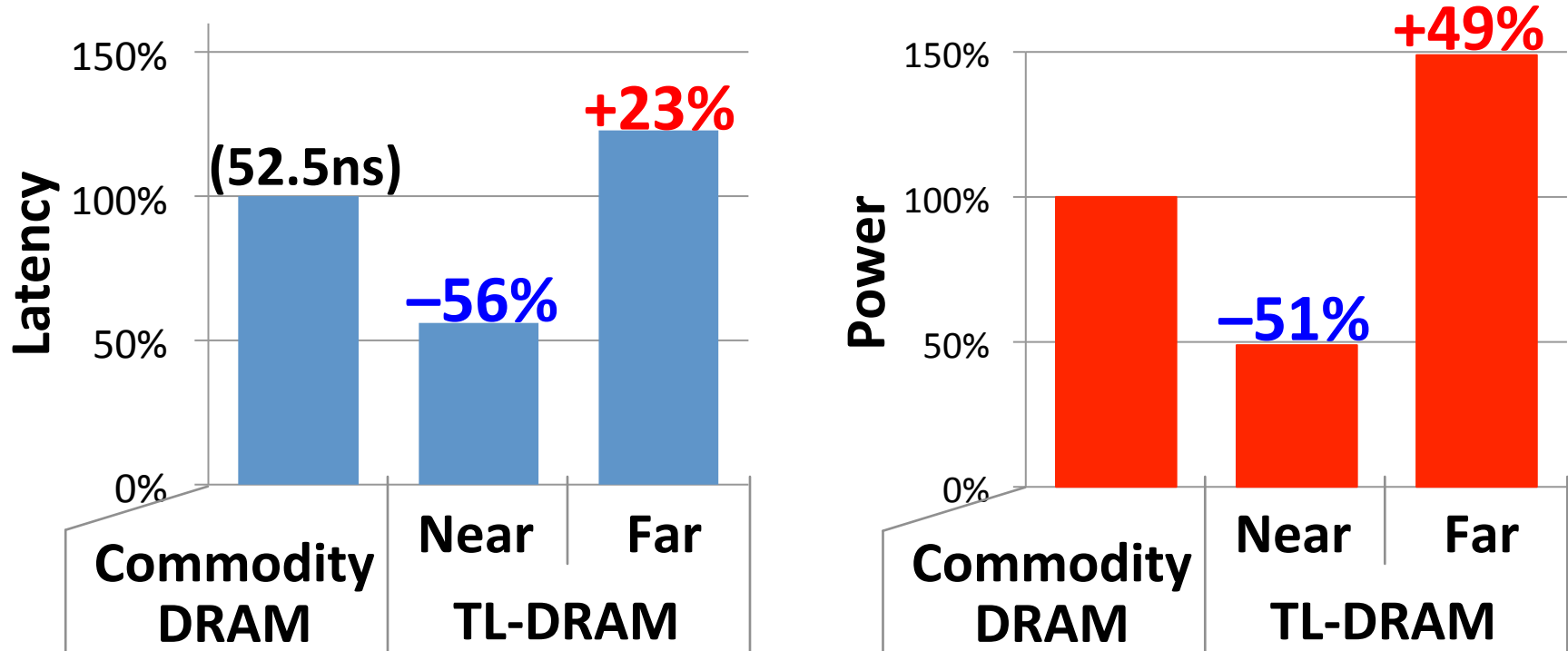
*Near Segment*

*Sense Amplifier*

# Latency, Power, and Area Evaluation

- **Commodity DRAM:** 512 cells/bitline
- **TL-DRAM:** 512 cells/bitline
  - Near segment: 32 cells
  - Far segment: 480 cells
- **Latency Evaluation**
  - SPICE simulation using circuit-level DRAM model
- **Power and Area Evaluation**
  - DRAM area/power simulator from Rambus
  - DDR3 energy calculator from Micron
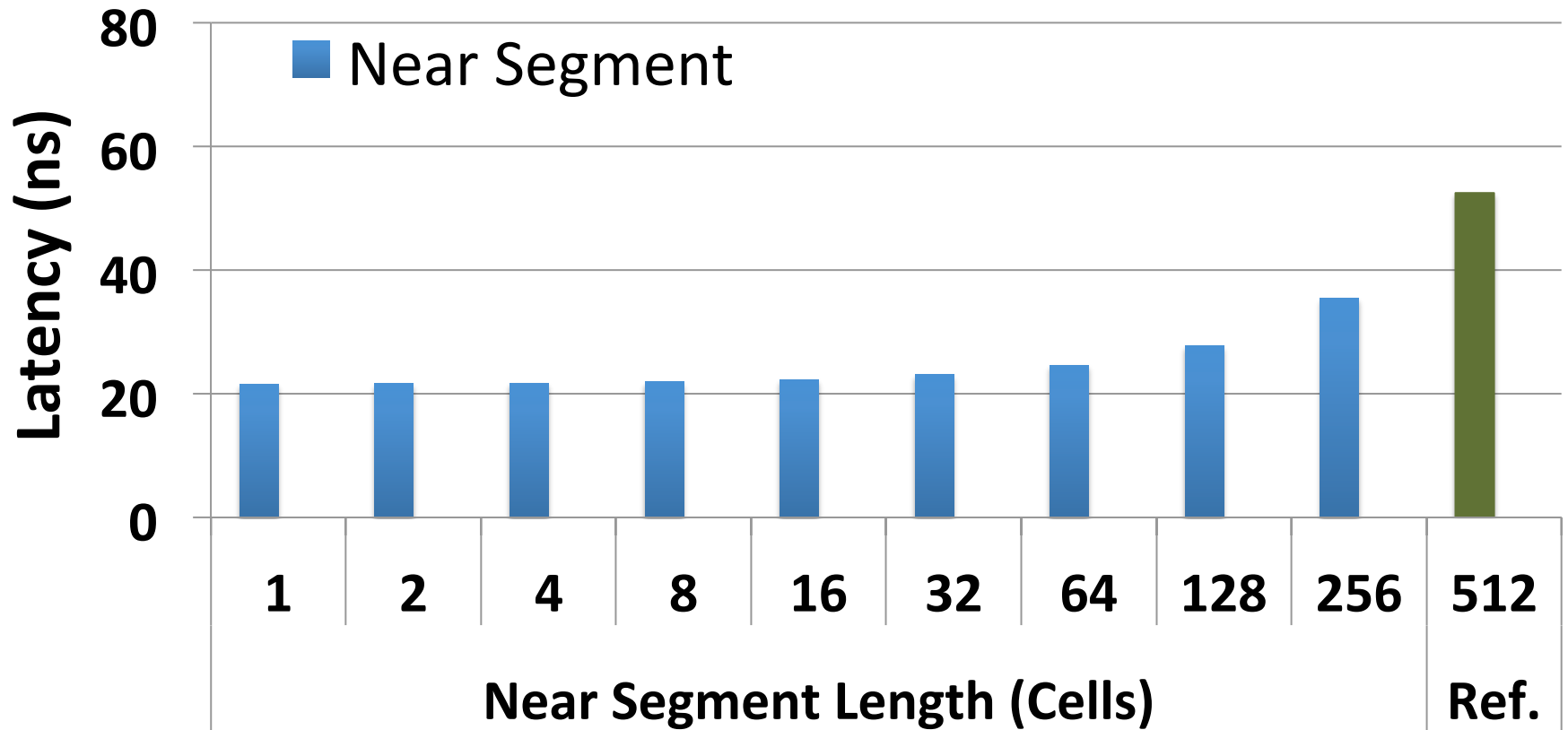
# Commodity DRAM vs. TL-DRAM

- **DRAM Latency** (tRC)
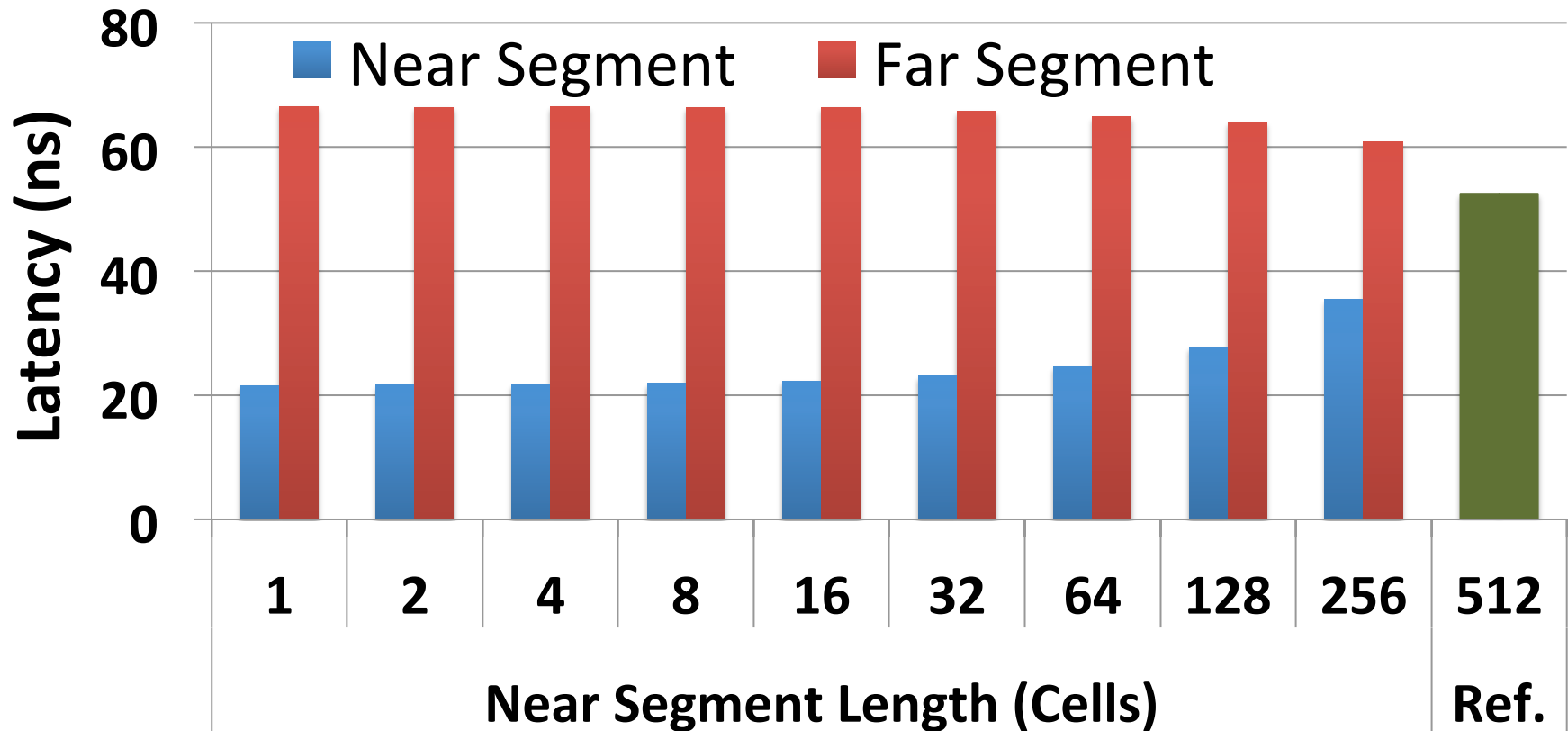- **DRAM Power**



- **DRAM Area Overhead**

**~3%**: mainly due to the isolation transistors

# Latency vs. Near Segment Length



*Longer near segment length leads to higher near segment latency*
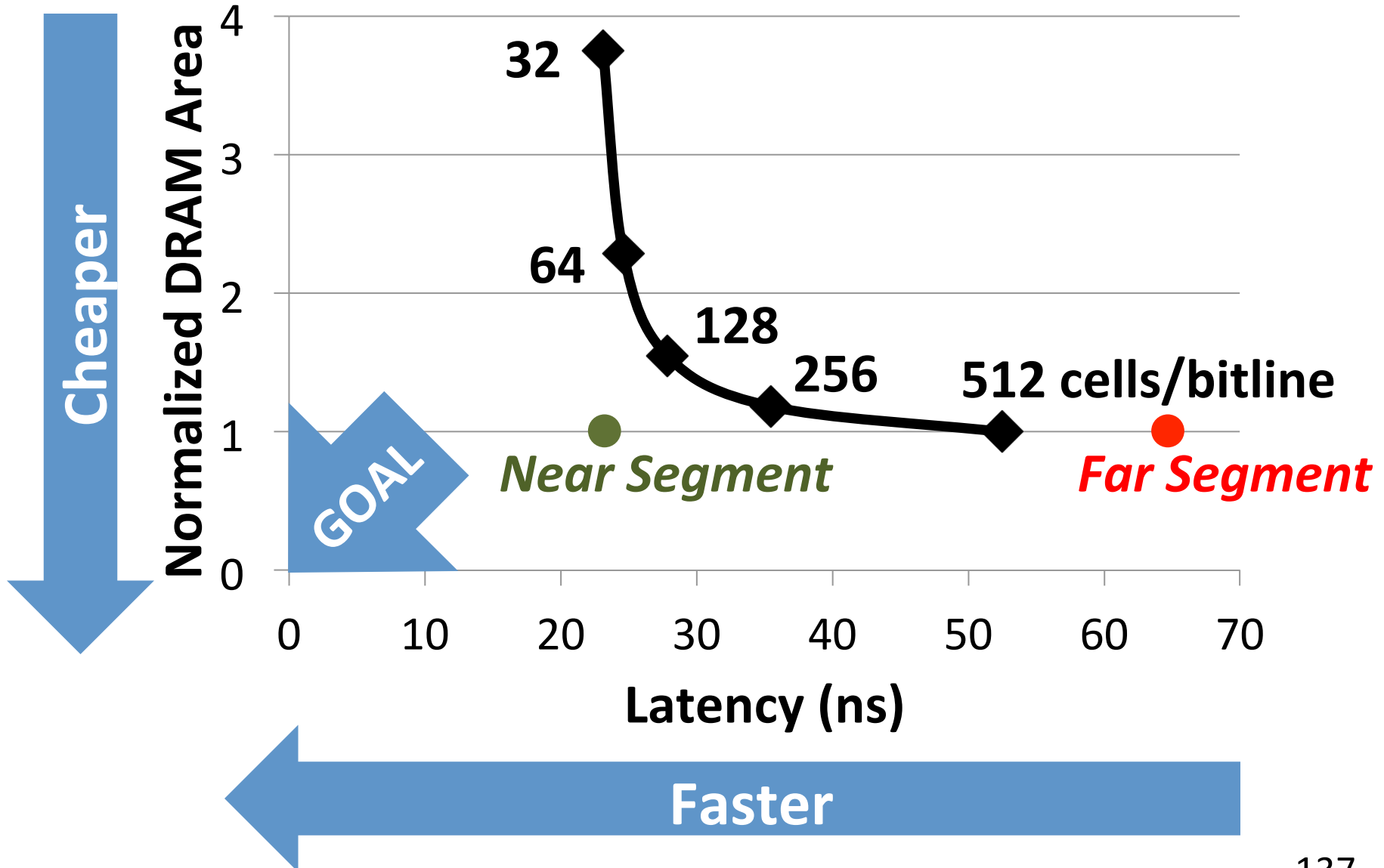
# Latency vs. Near Segment Length



Far segment latency is higher than commodity DRAM latency

# Trade-Off: Area (Die-Area) vs. Latency
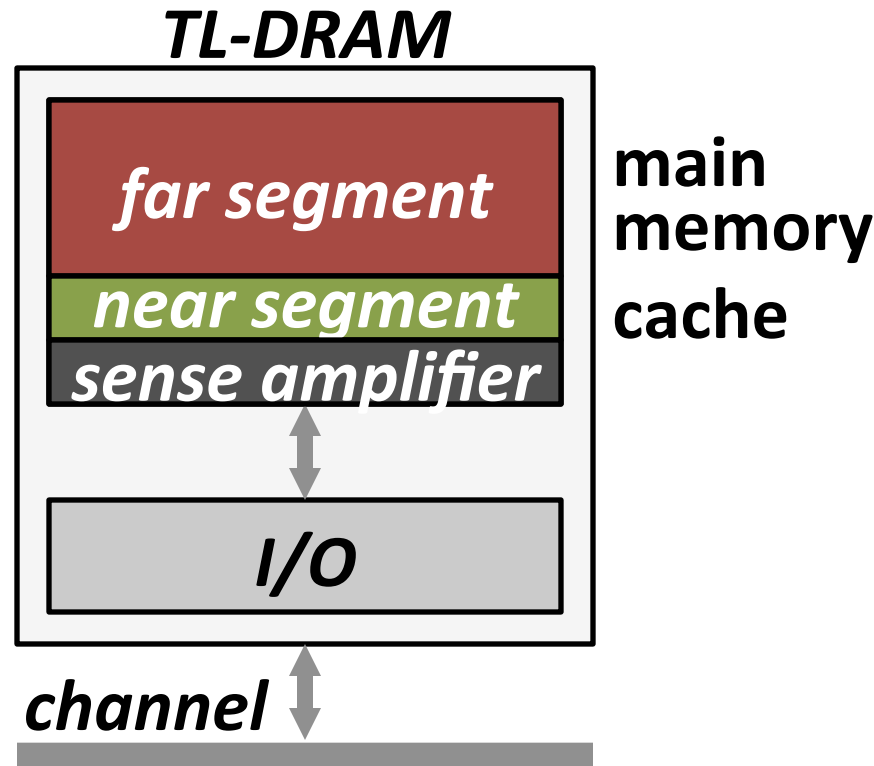
# Leveraging Tiered-Latency DRAM

- TL-DRAM is a **substrate** that can be leveraged by the hardware and/or software

- Many potential uses
  1. Use near segment as hardware-managed *inclusive* cache to far segment
  2. Use near segment as hardware-managed *exclusive* cache to far segment
  3. Profile-based page mapping by operating system
  4. Simply replace DRAM with TL-DRAM

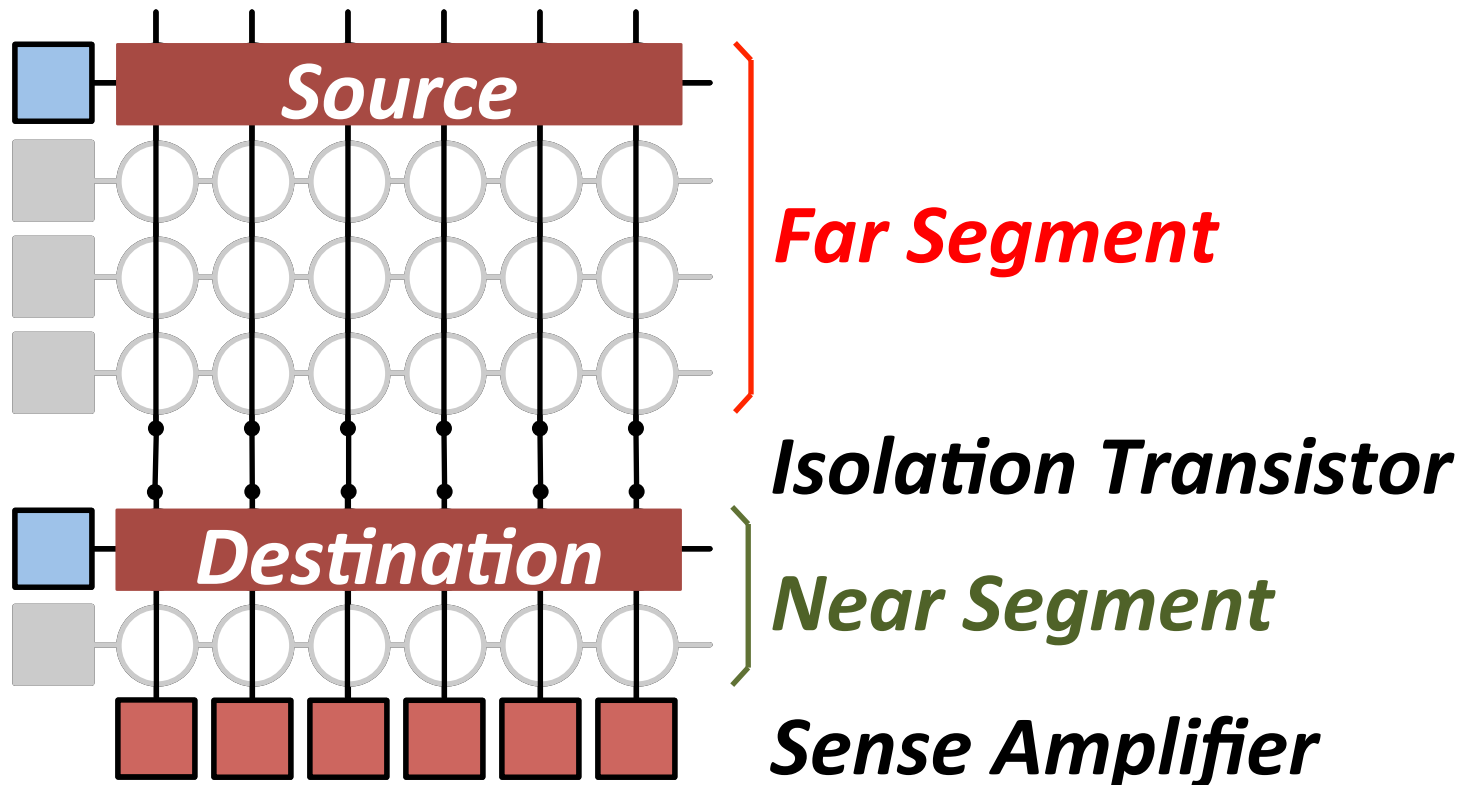# Near Segment as Hardware-Managed Cache

**TL-DRAM**

far segment — main memory

near segment — cache

sense amplifier

I/O

channel

- **Challenge 1:** How to efficiently migrate a row between segments?

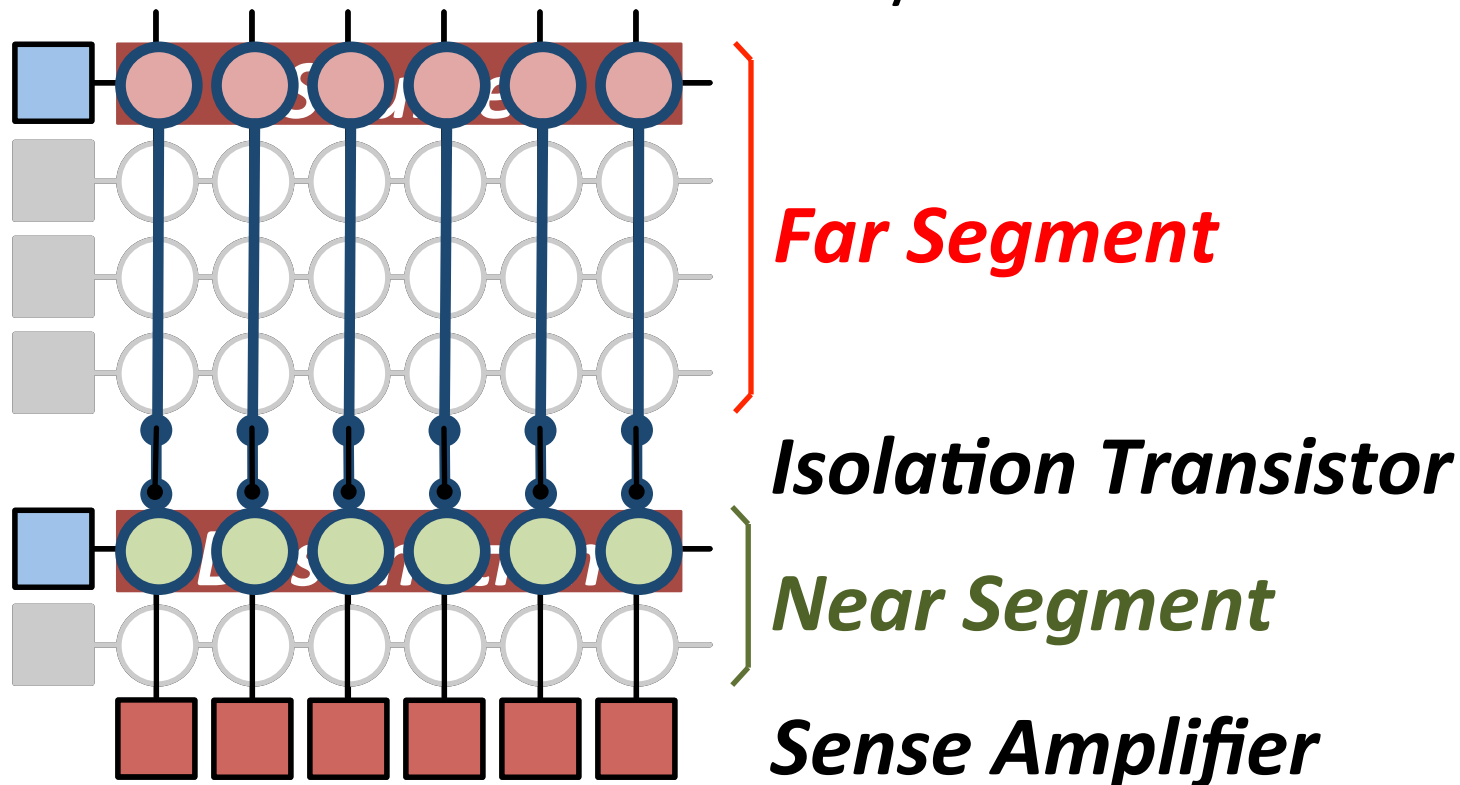- **Challenge 2:** How to efficiently manage the cache?

# Inter-Segment Migration

- **Goal:** Migrate source row into destination row
- **Naïve way:** Memory controller reads the source row *byte by byte* and writes to destination row *byte by byte*
  → *High latency*



**Far Segment**

**Isolation Transistor**

**Near Segment**

**Sense Amplifier**

# Inter-Segment Migration

- **Our way:**
  - Source and destination cells *share bitlines*
  - Transfer data from source to destination across ***shared bitlines*** concurrently



*Far Segment*

*Isolation Transistor*

*Near Segment*

*Sense Amplifier*

# Inter-Segment Migration

- **Our way:**
  - Source and destination cells *share bitlines*
  - Transfer data from so...
    *shared bitlines* concur...

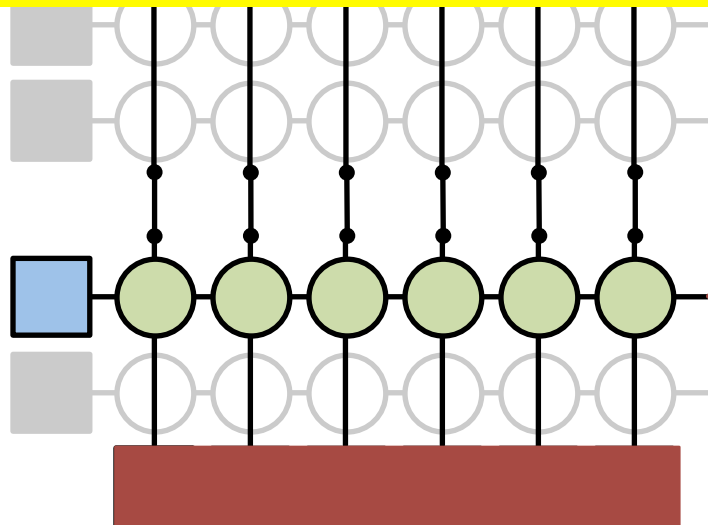**Step 1:** Activate source row

**Migration is overlapped with source row access**

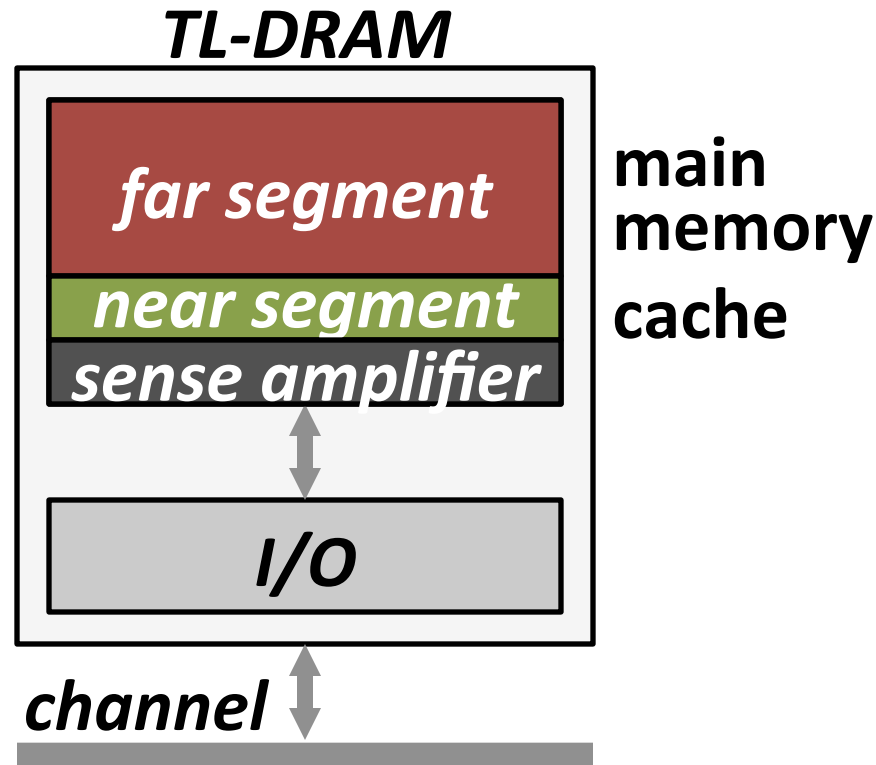**Additional ~4ns over row access latency**

**Step 2:** Activate destination row to connect cell and bitline

*Iso...ll Transistor*

*Near Segment*

*Sense Amplifier*

# Near Segment as Hardware-Managed Cache

**TL-DRAM**



- **Challenge 1:** How to efficiently migrate a row between segments?
- **Challenge 2:** How to efficiently manage the cache?

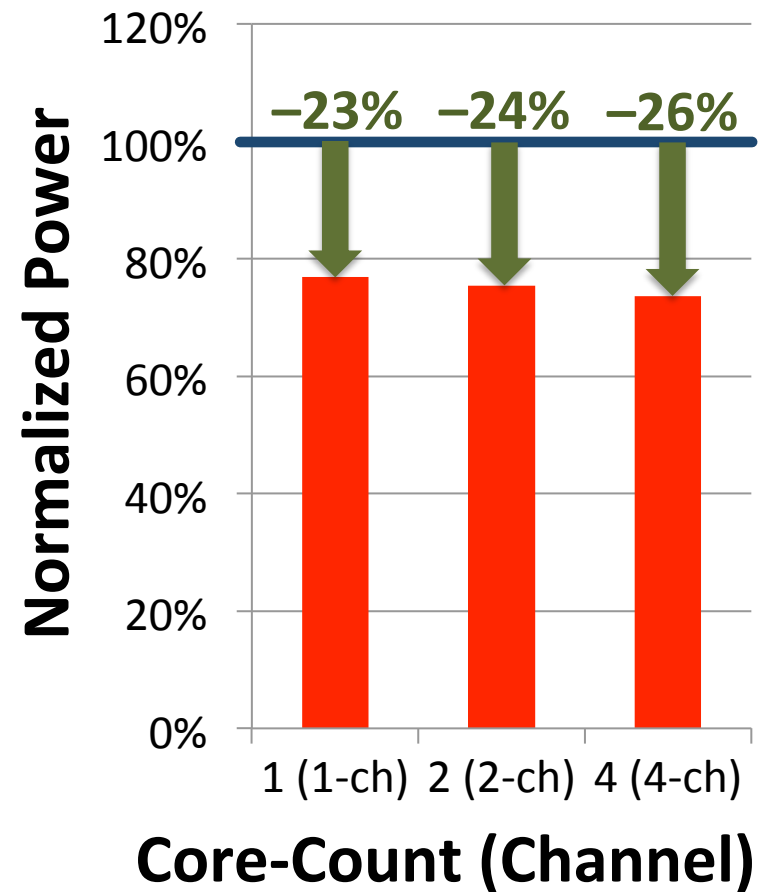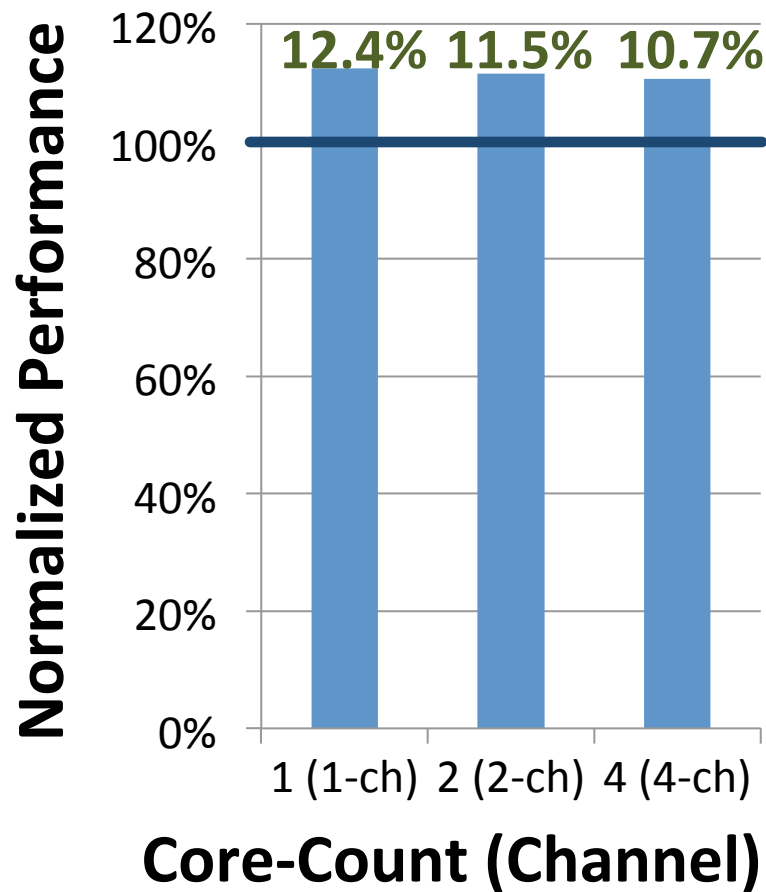# Evaluation Methodology

- **System simulator**
  - CPU: Instruction-trace-based x86 simulator
  - Memory: Cycle-accurate DDR3 DRAM simulator

- **Workloads**
  - 32 Benchmarks from TPC, STREAM, SPEC CPU2006

- **Performance Metrics**
  - Single-core: Instructions-Per-Cycle
  - Multi-core: Weighted speedup

# Configurations

- **System configuration**
  - CPU: 5.3GHz
  - LLC: 512kB private per core
  - **Memory: DDR3-1066**
    - 1-2 channel, 1 rank/channel
    - 8 banks, 32 subarrays/bank, **512 cells/bitline**
    - Row-interleaved mapping & closed-row policy

- **TL-DRAM configuration**
  - Total bitline length: **512 cells/bitline**
  - Near segment length: 1-256 cells
  - Hardware-managed inclusive cache: near segment

# Performance & Power Consumption



*Using near segment as a cache improves performance and reduces power consumption*

146

# Single-Core: Varying Near Segment Length



*By adjusting the near segment length, we can trade off cache capacity for cache latency*

# Other Mechanisms & Results

- **More mechanisms** for leveraging TL-DRAM
  - Hardware-managed *exclusive* caching mechanism
  - Profile-based page mapping to near segment
  - TL-DRAM improves performance and reduces power consumption with other mechanisms
- **More than two tiers**
  - Latency evaluation for three-tier TL-DRAM
- **Detailed circuit evaluation** for DRAM latency and power consumption
  - Examination of tRC and tRCD
- **Implementation details** and **storage cost analysis** in memory controller

# Summary of TL-DRAM

- **Problem: DRAM latency is a critical performance bottleneck**
- **Our Goal**: Reduce DRAM latency with low area cost
- **Observation**: Long bitlines in DRAM are the dominant source of DRAM latency
- **Key Idea: Divide long bitlines into two shorter segments**
  - **Fast and slow segments**
- **Tiered-latency DRAM: Enables latency heterogeneity in DRAM**
  - **Can leverage this in many ways to improve performance and reduce power consumption**
- **Results**: When the fast segment is used as a cache to the slow segment → Significant performance improvement (>12%) and power reduction (>23%) at low area cost (3%)

# New DRAM Architectures

- RAIDR: Reducing Refresh Impact
- TL-DRAM: Reducing DRAM Latency
- SALP: Reducing Bank Conflict Impact
- RowClone: Fast Bulk Data Copy and Initialization

# Subarray-Level Parallelism: Reducing Bank Conflict Impact

Yoongu Kim, Vivek Seshadri, Donghyuk Lee, Jamie Liu, and Onur Mutlu,
**"A Case for Exploiting Subarray-Level Parallelism (SALP) in DRAM"**
*Proceedings of the 39th International Symposium on Computer Architecture (**ISCA**),*
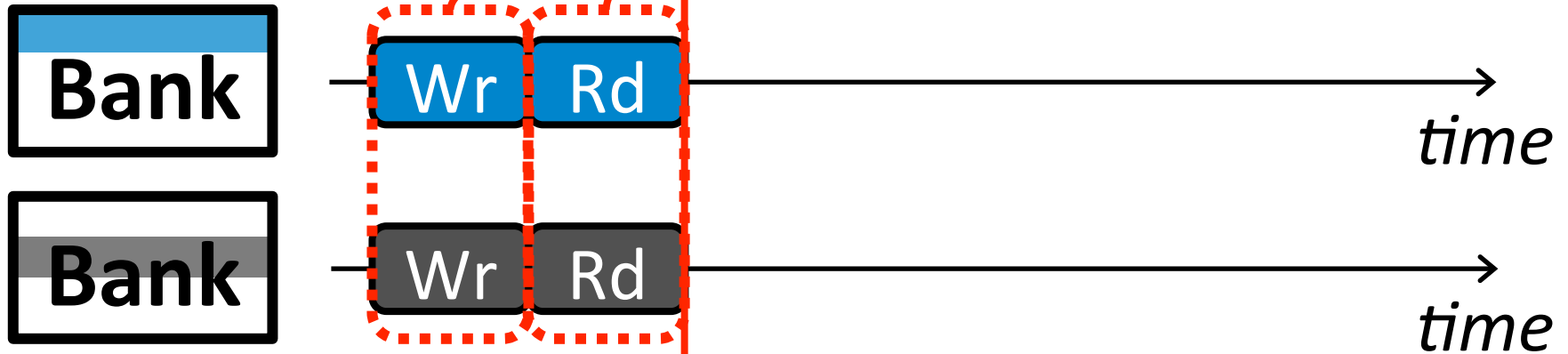Portland, OR, June 2012. Slides (pptx)

# The Memory Bank Conflict Problem

- Two requests to the same bank are serviced serially

- Problem: Costly in terms of performance and power

- Goal: We would like to reduce bank conflicts without increasing the number of banks (at low cost)

- Idea: Exploit the internal sub-array structure of a DRAM bank to parallelize bank conflicts
  - By reducing global sharing of hardware between sub-arrays

- Kim, Seshadri, Lee, Liu, Mutlu, "A Case for Exploiting Subarray-Level Parallelism in DRAM," ISCA 2012.

**SAFARI**

# The Problem with Memory Bank Conflicts
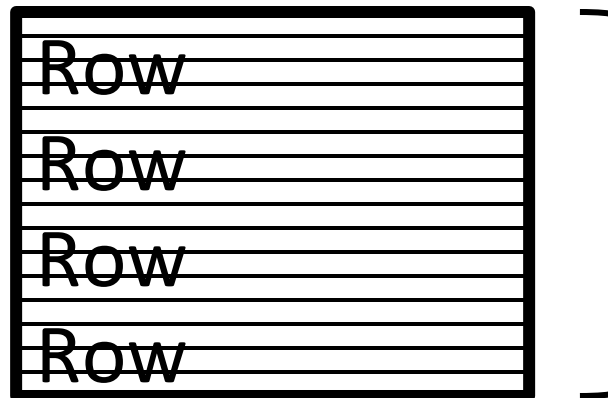
- **Two Banks**



- **One Bank**



153

# Goal

- **Goal:** *Mitigate the detrimental effects of bank conflicts in a cost-effective manner*

- **Naïve solution:** Add more banks
  - Very expensive

- **Cost-effective solution:** Approximate the benefits of more banks without adding more banks

# Key Observation #1

A DRAM bank is divided into *subarrays*

*Logical Bank*  *Physical Bank*

Row
Row
Row
Row

$32k\ rows$

Subarray$_{64}$

Subarray$_{1}$

Row-Buffer

Global Row-Buf
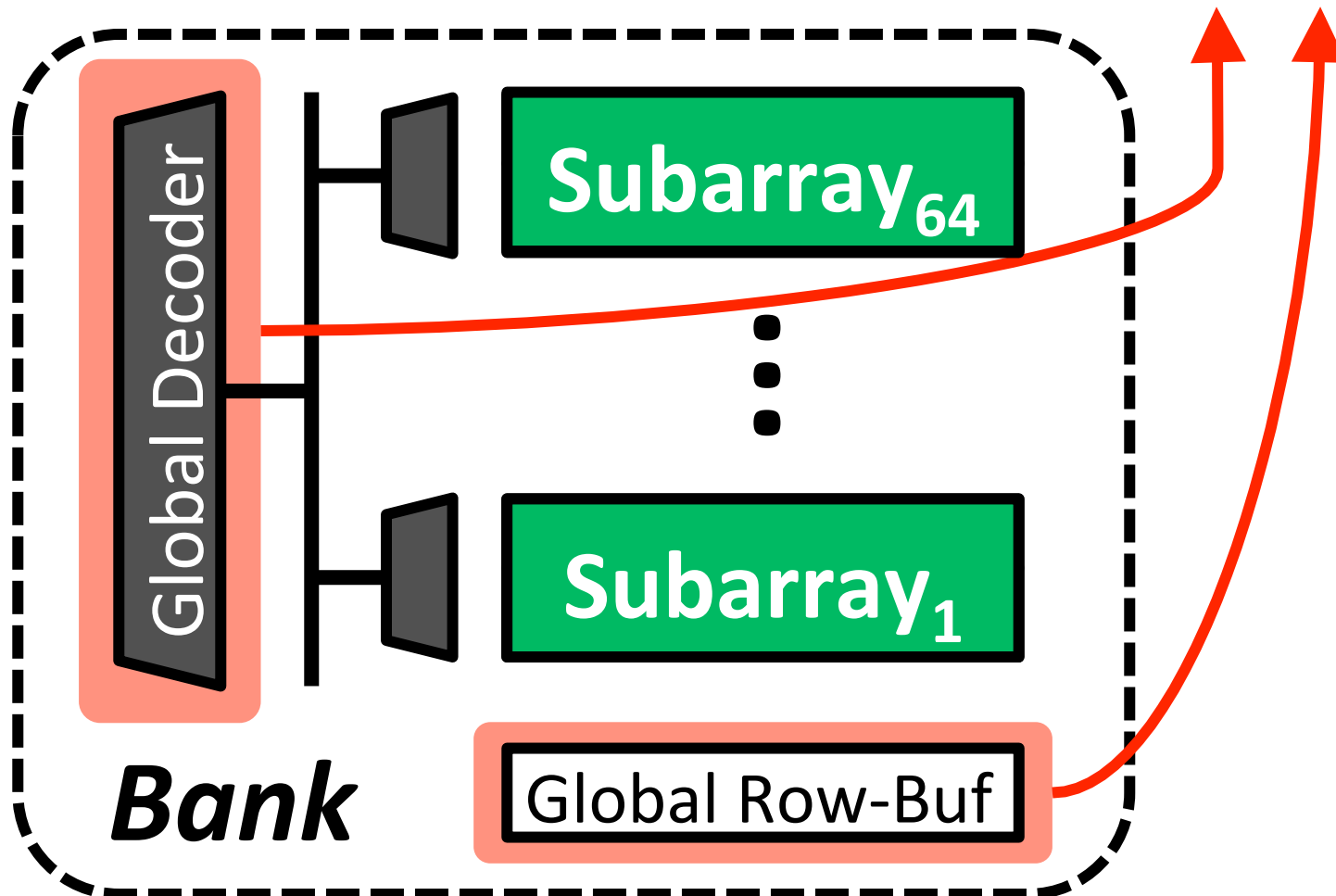
A *single* row-buffer cannot drive *all* rows

Many *local row-buffers*, one at each *subarray*

# Key Observation #2

Each subarray is mostly independent...
– except occasionally sharing *global structures*

# Key Idea: Reduce Sharing of Globals

1. Parallel access to subarrays



Global Decoder

Local Row-Buf

Local Row-Buf

*Bank*

Global Row-Buf

2. Utilize multiple local row-buffers

# Overview of Our Mechanism

Subarray₆₄

⋮

Subarray₁

Global Row-Buf

*1. Parallelize*

Reg Reg Reg Reg

*2. Utilize multiple*
*To same bank...*
*local row-buffers*
*but diff. subarrays*

# Challenges: Global Structures

## 1. Global Address Latch

## 2. Global Bitlines

# Challenge #1. Global Address Latch

# Solution #1. Subarray Address Latch



ACTIVATED

ACTIVATED

$V_{DD}$

$V_{DD}$

Local row-buffer

Local row-buffer

Global row-buffer

*Global latch → local latches*

# Challenges: Global Structures

## 1. Global Address Latch

- Problem: Only *one* raised wordline
- Solution: **Subarray Address Latch**

## 2. Global Bitlines

# Challenge #2. Global Bitlines

# Solution #2. Designated-Bit Latch



**Global bitlines**

Wire

*Local row-buffer*

**D**

*Switch*

*Local row-buffer*

**D**

*Switch*

**READ**

*Global row-buffer*

*Selectively connect local to global*

# **Challenges: Global Structures**

## **1. Global Address Latch**

- Problem: Only *one* raised wordline

- Solution: **Subarray Address Latch**

## **2. Global Bitlines**

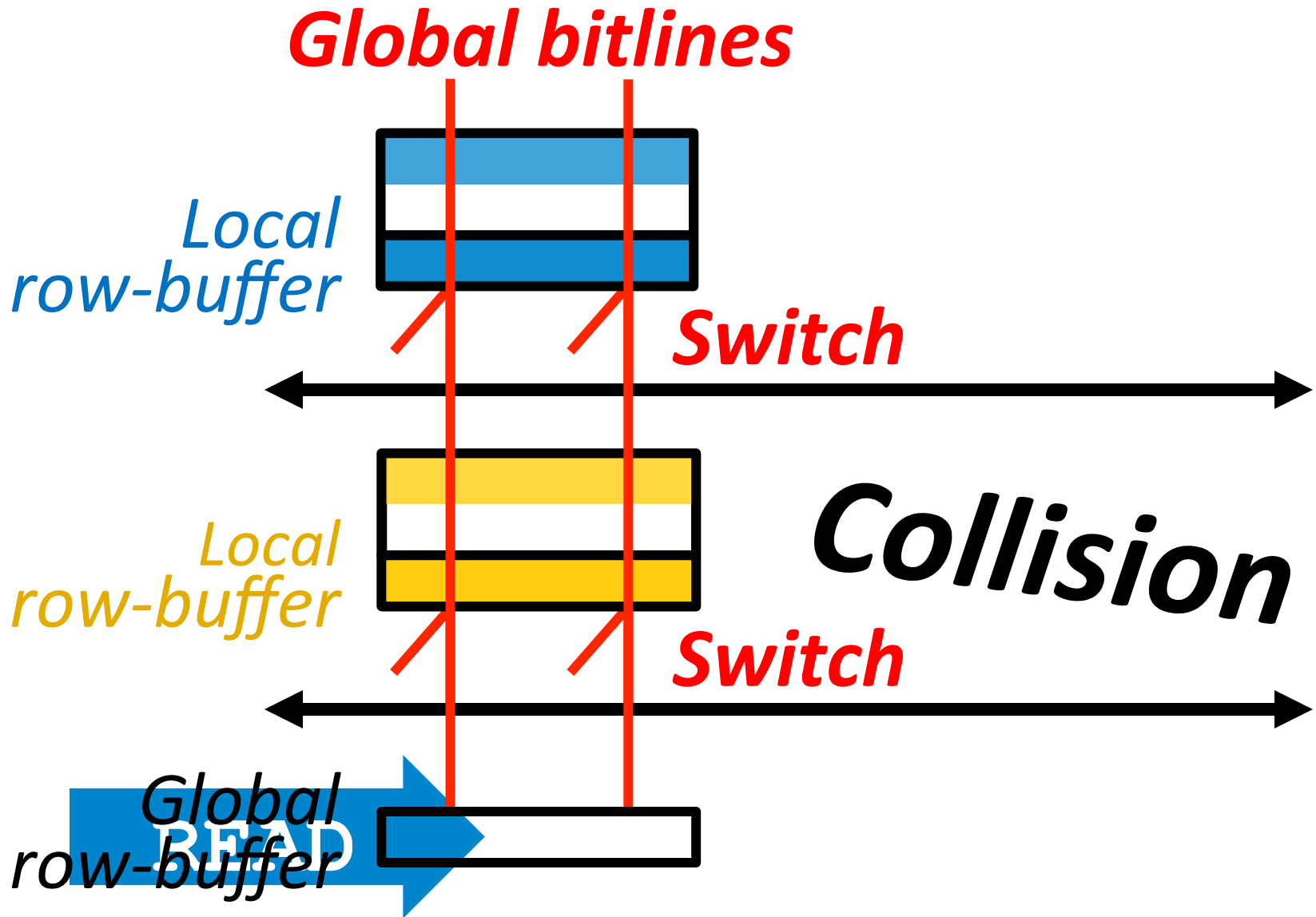- Problem: Collision during access

- Solution: **Designated-Bit Latch**

**MASA (Multitude of Activated Subarrays)**

# **MASA: Advantages**

- Baseline (Subarray-Oblivious)



- **MASA**

# MASA: Overhead

- **DRAM Die Size**: Only **0.15%** increase
  - Subarray Address Latches
  - Designated-Bit Latches & Wire
- **DRAM Static Energy:** Small increase
  - **0.56mW** for each activated subarray
  - *But saves dynamic energy*
- **Controller**: Small additional storage
  - Keep track of subarray status (< **256B**)
  - Keep track of new timing constraints

# Cheaper Mechanisms

*Latches*

1. Serialization
2. Wr-Penalty
3. Thrashing

MASA

SALP-2

SALP-1

# System Configuration

- **System Configuration**
  - CPU: 5.3GHz, 128 ROB, 8 MSHR
  - LLC: 512kB per-core slice

- **Memory Configuration**
  - DDR3-1066
  - *(default)* **1 channel, 1 rank, 8 banks, 8 subarrays-per-bank**
  - *(sensitivity)* 1-8 chans, 1-8 ranks, 8-64 banks, 1-128 subarrays

- **Mapping & Row-Policy**
  - *(default)* **Line-interleaved & Closed-row**
  - *(sensitivity)* Row-interleaved & Open-row

- **DRAM Controller Configuration**
  - 64-/64-entry read/write queues per-channel
  - FR-FCFS, batch scheduling for writes

# SALP: Single-core Results



**MASA** *achieves most of the benefit of having more banks ("Ideal")*

170

# SALP: Single-Core Results



SALP-1, SALP-2, MASA improve performance at low cost

# Subarray-Level Parallelism: Results



*MASA increases energy-efficiency*

# New DRAM Architectures

- RAIDR: Reducing Refresh Impact
- TL-DRAM: Reducing DRAM Latency
- SALP: Reducing Bank Conflict Impact
- RowClone: Fast Bulk Data Copy and Initialization

SAFARI

# RowClone: Fast Bulk Data Copy and Initialization

Vivek Seshadri, Yoongu Kim, Chris Fallin, Donghyuk Lee, Rachata Ausavarungnirun,
Gennady Pekhimenko, Yixin Luo, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry,
**"RowClone: Fast and Efficient In-DRAM Copy and Initialization of Bulk Data"**
*CMU Computer Science Technical Report*, CMU-CS-13-108, Carnegie Mellon University, April 2013.

# Today's Memory: Bulk Data Copy



3) Cache pollution

1) High latency

2) High bandwidth utilization

4) Unwanted data movement

Memory

CPU    L1    L2    L3    MC

# Future: RowClone (In-Memory Copy)



3) No cache pollution

1) Low latency

**Memory**

**CPU** — **L1** **L2** **L3** — **MC**

2) Low bandwidth utilization

4) No unwanted data movement

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

# DRAM operation (load one byte)

4 Kbits

1. Activate row

2. Transfer row

DRAM array

Row Buffer (4 Kbits)

3. Transfer byte onto bus

Data pins (8 bits)

Memory Bus

# RowClone: in-DRAM Row Copy (and Initialization)



4 Kbits

1. Activate row A

3. Activate row B

DRAM array

2. Transfer row

4. Transfer row

Row Buffer (4 Kbits)

Data pins (8 bits)

Memory Bus

# Our Approach: Key Idea

- DRAM banks contain
  1. Mutiple rows of DRAM cells – row = 8KB
  2. A row buffer shared by the DRAM rows

- Large scale copy
  1. Copy data from source row to row buffer
  2. Copy data from row buffer to destination row

# DRAM Subarray Microarchitecture

DRAM Row
(share wordline)
(~8Kb)

Sense
Amplifiers
(row buffer)

DRAM Cell

wordline

# DRAM Operation

Raise wordline

| ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? | src |
|---|---|---|---|---|---|---|---|---|---|---|---|-----|
|   |   |   |   |   |   |   |   |   |   |   |   |     |
| 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | dst |
|   |   |   |   |   |   |   |   |   |   |   |   |     |

Sense Amplifiers (row buffer)

| - | + | - | - | + | + | - | - | - | + | + | - |
|---|---|---|---|---|---|---|---|---|---|---|---|

Activate (src) → Precharge

181

# RowClone: Intra-subarray Copy



|   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | **1** | **0** | **0** | **1** | **1** | **0** | **0** | **0** | **1** | **1** | **0** | src
| **0** | **1** | **0** | **0** | **1** | **1** | **0** | **0** | **0** | **1** | **1** | **0** | dst

Sense
Amplifiers
(row buffer)

| **0** | **1** | **0** | **0** | **1** | **1** | **0** | **0** | **0** | **1** | **1** | **0** |

Activate (src) ⟶ Deactivate (our proposal) ⟶ Activate (dst)

# RowClone: Inter-bank Copy



src

dst

Read

Write

I/O Bus

Transfer
(our proposal)

# RowClone: Inter-subarray Copy



dst

src

temp

I/O Bus

1. Transfer (src to temp)
2. Transfer (temp to dst)

# Fast Row Initialization



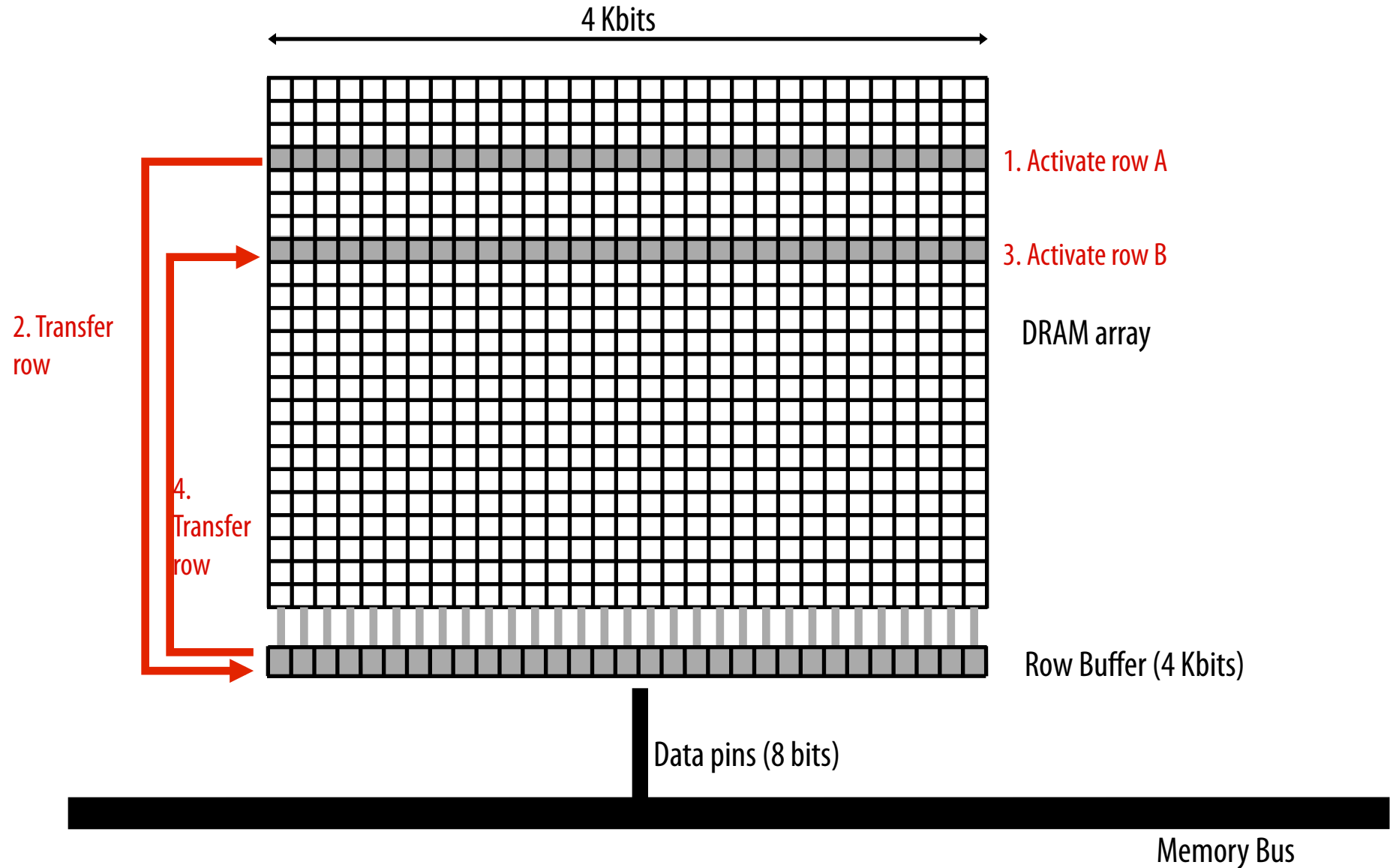**Fix a row at Zero**
(0.5% loss in capacity)

# RowClone: Latency and Energy Savings



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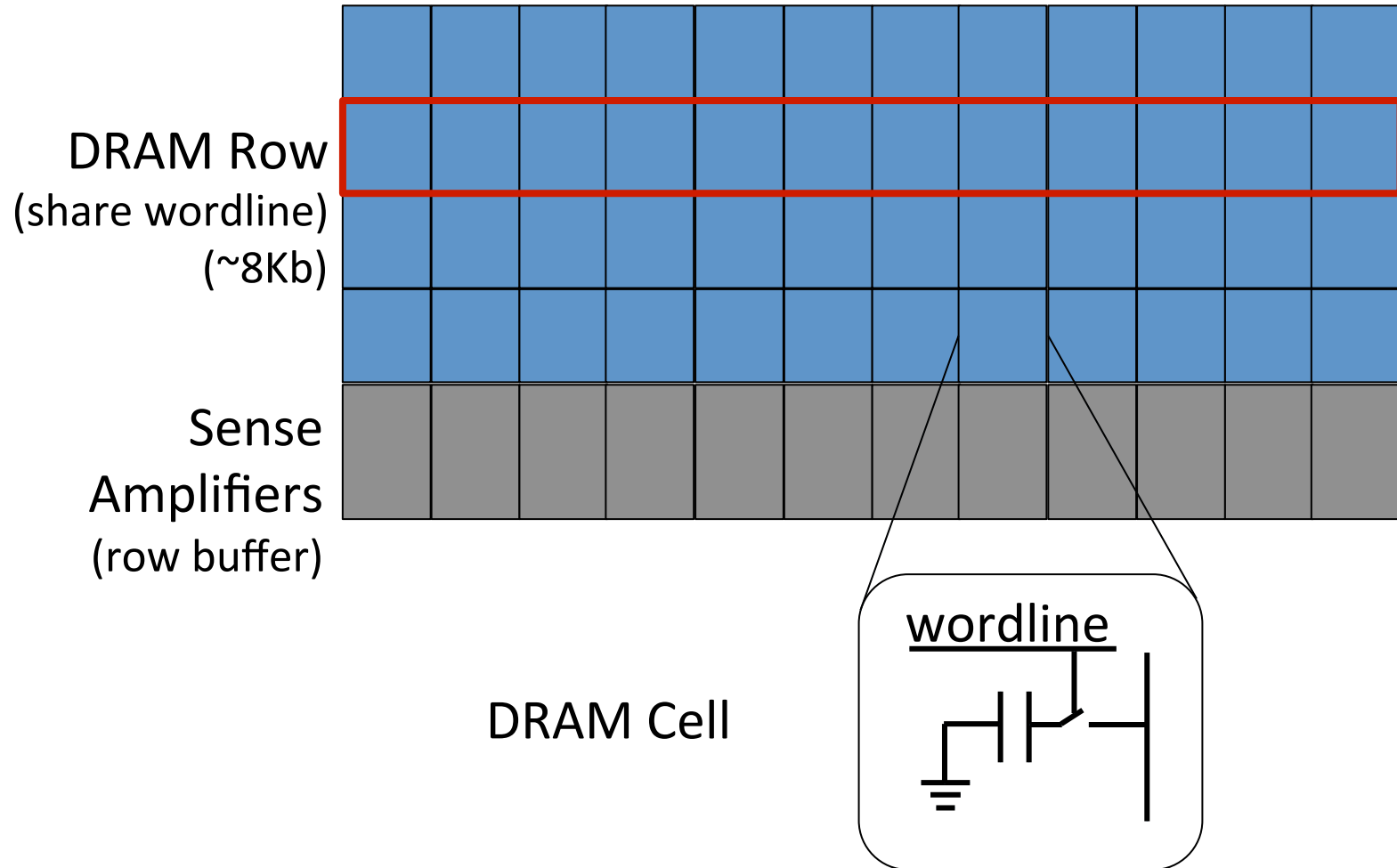
# Summary

- Major problems with DRAM scaling and design: high refresh rate, high latency, low parallelism, bulk data movement

- Four new DRAM designs
  - RAIDR: Reduces refresh impact
  - TL-DRAM: Reduces DRAM latency at low cost
  - SALP: Improves DRAM parallelism
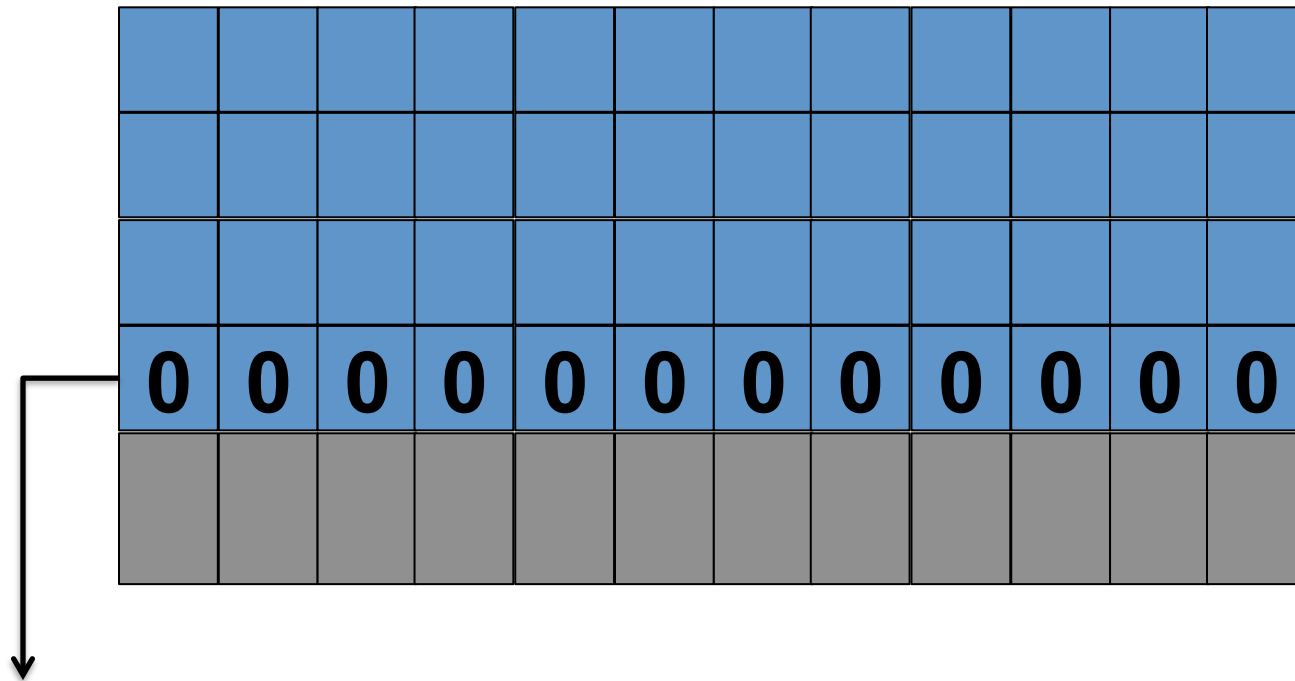  - RowClone: Reduces energy and performance impact of bulk data copy

- All four designs
  - Improve both performance and energy consumption
  - Are low cost (low DRAM area overhead)
  - Enable new degrees of freedom to software & controllers

- Rethinking DRAM interface and design essential for scaling
  - Co-design DRAM with the rest of the system

**SAFARI**

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

**Carnegie Mellon**

SAFARI

# Additional Material

# Three Papers

- 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"**
*<u>Intel Technology Journal</u>* (**ITJ**) *Special Issue on Memory Resiliency*, Vol. 17, No. 1, May 2013.

- Howard David, Chris Fallin, Eugene Gorbatov, Ulf R. Hanebutte, and <u>Onur Mutlu</u>,
**"Memory Power Management via Dynamic Voltage/Frequency Scaling"**
*Proceedings of the <u>8th International Conference on Autonomic Computing</u>* (**ICAC**), Karlsruhe, Germany, June 2011. <u>Slides (pptx)</u> <u>(pdf)</u>

- Jamie Liu, Ben Jaiyen, Yoongu Kim, Chris Wilkerson, and <u>Onur Mutlu</u>,
**"An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms"**
*Proceedings of the <u>40th International Symposium on Computer Architecture</u>* (**ISCA**), Tel-Aviv, Israel, June 2013. <u>Slides (pptx)</u> <u>Slides (pdf)</u>

# Aside: Scaling Flash Memory [Cai+, ICCD'12]

- NAND flash memory has low endurance: a flash cell dies after 3k P/E cycles vs. 50k desired → Major scaling challenge for flash memory

- Flash error rate increases exponentially over flash lifetime

- Problem: Stronger error correction codes (ECC) are ineffective and undesirable for improving flash lifetime due to
  - diminishing returns on lifetime with increased correction strength
  - prohibitively high power, area, latency overheads

- Our Goal: Develop techniques to tolerate high error rates w/o strong ECC

- Observation: Retention errors are the dominant errors in MLC NAND flash
  - flash cell loses charge over time; retention errors increase as cell gets worn out

- Solution: Flash Correct-and-Refresh (FCR)
  - Periodically read, correct, and reprogram (in place) or remap each flash page before it accumulates more errors than can be corrected by simple ECC
  - Adapt "refresh" rate to the severity of retention errors (i.e., # of P/E cycles)

- Results: FCR improves flash memory lifetime by 46X with no hardware changes and low energy overhead; outperforms strong ECCs

# DRAM Experiments: Summary (I)

- DRAM requires periodic refresh to avoid data loss
  - Refresh wastes energy, reduces performance, limits DRAM density scaling
- Many past works observed that different DRAM cells can retain data for different times without being refreshed; proposed reducing refresh rate for strong DRAM cells
  - Problem: These techniques require an accurate profile of the retention time of all DRAM cells
- Our goal: To analyze the retention time behavior of DRAM cells in modern DRAM devices to aid the collection of accurate profile information
- Our experiments: We characterize 248 modern commodity DDR3 DRAM chips from 5 manufacturers using an FPGA based testing platform
- Two Key Issues:

  1. Data Pattern Dependence: A cell's retention time is heavily dependent on data values stored in itself and nearby cells, which cannot easily be controlled.

  2. Variable Retention Time: Retention time of some cells change unpredictably from high to low at large timescales.

# DRAM Experiments: Summary (II)

- **Key findings on Data Pattern Dependence**
  - There is no observed single data pattern that elicits the lowest retention times for a DRAM device → very hard to find this pattern
  - DPD varies between devices due to variation in DRAM array circuit design between manufacturers
  - DPD of retention time gets worse as DRAM scales to smaller feature sizes

- **Key findings on Variable Retention Time**
  - VRT is common in modern DRAM cells that are weak
  - The timescale at which VRT occurs is very large (e.g., a cell can stay in high retention time state for a day or longer) → finding minimum retention time can take very long

- Future work on retention time profiling must address these issues

**SAFARI**