

MEMORY SYSTEMS

Mahdi Nazm Bojnordi

Assistant Professor

School of Computing

University of Utah

Overview

- Upcoming deadline
 - ▣ Feb. 27th: homework assignment will be posted
- Group projects

	Subject	Presentation
Group 1	Reducing cache energy via skipping common values	April 17
Group 2	Integrating high bandwidth memory with high capacity DIMMs for data-centric workloads	April 17
Group 3	A near-threshold spiking neural network accelerator	April 17
Group 4	Improving energy efficiency of low power micro-controllers used in various IoT nodes	April 19
Group 5	In-situ hardware accelerator for clustering applications	April 19
Group 6	Improving performance of time-based data encoding	April 19

Overview

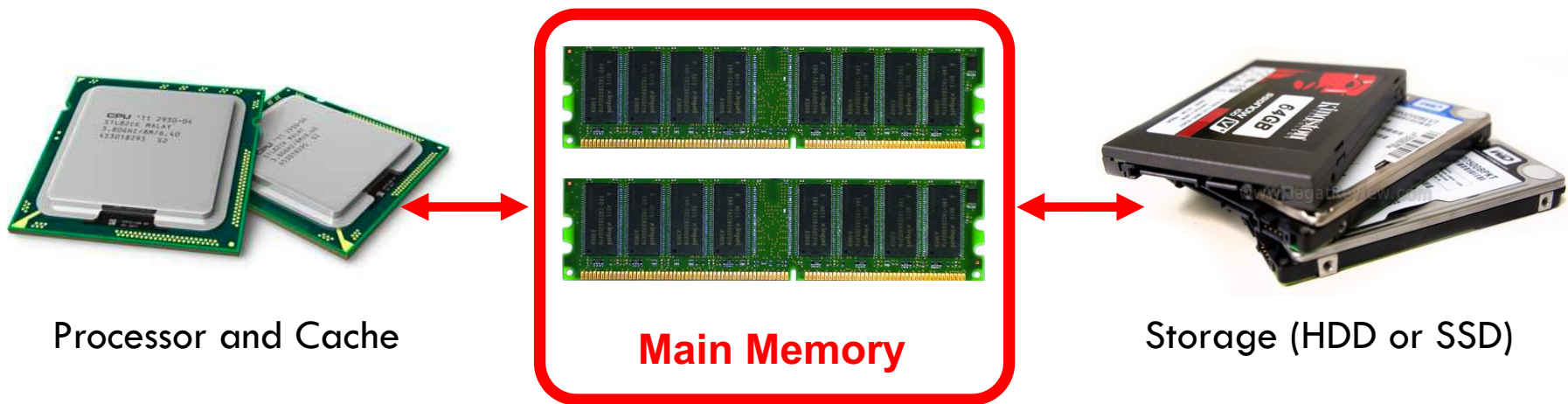
- Student paper presentation
 - ▣ Prepare for **exactly** 20m talk followed by 5m Q&A

Presenter Names	Date
Kohl, Meher, Shirley	March 29
Karl, Anirban, Chandrasekhar	April 3
Suryanarayanan, Tim, Arjun	April 5
Pranav, Goverdhan, Yomi	April 10
Munzer, Manikanth, Amandeep	April 12

- This lecture
 - ▣ Main memory systems
 - ▣ DRAM architecture basics

Main Memory System

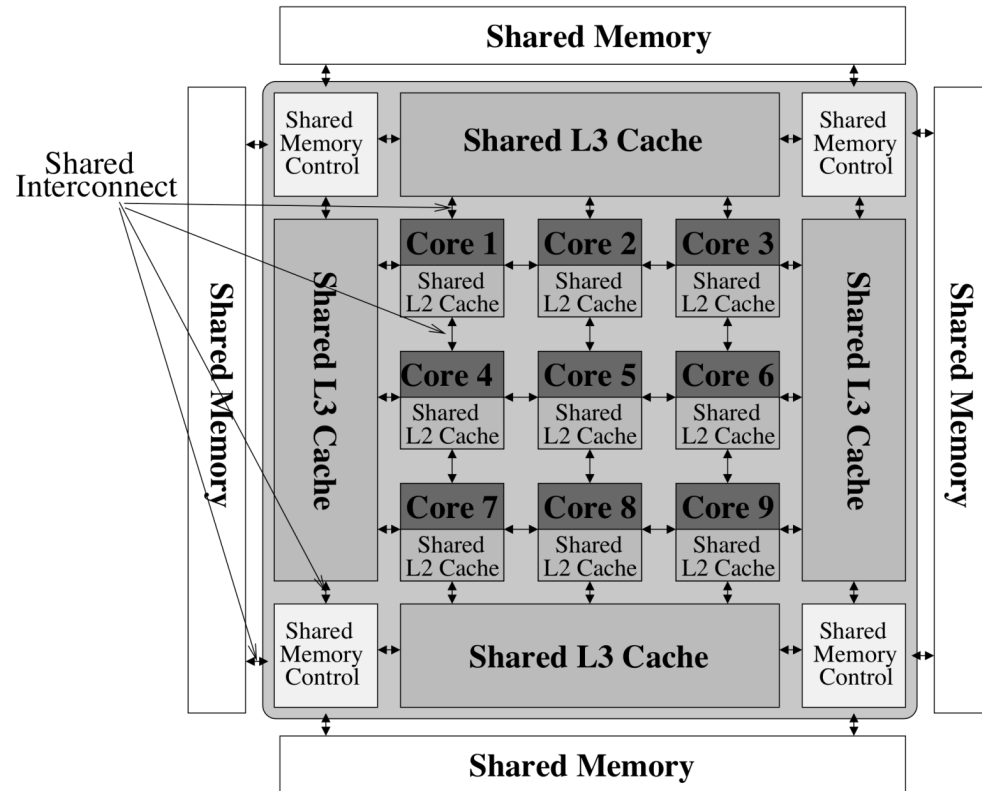
- A critical component of all computing systems
 - ▣ server, mobile, embedded, desktop, sensor



- Must scale to maintain performance growth
 - ▣ size, technology, efficiency, cost, and control algorithms

Why Main Memory is Important?

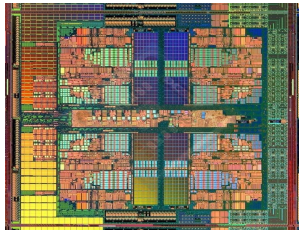
- Shared resource
 - ▣ Multiple applications running on different processor cores
 - ▣ Different objectives and requirements
 - ▣ Highly contented resource



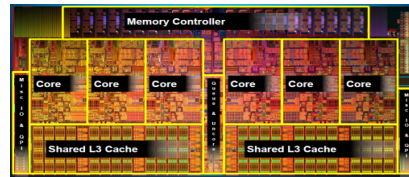
Complex control policies and microarchitectures are required.

Scalability Challenges

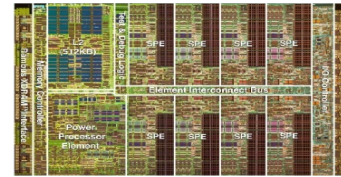
- Increasing need for memory capacity, bandwidth, and quality-of-service maintenance
 - ▣ increasing number of cores (multicores)
 - ▣ increasing demand for data (big data processing)
 - ▣ cloud computing, GPUs, mobile (consolidation)



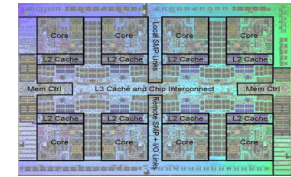
AMD Barcelona



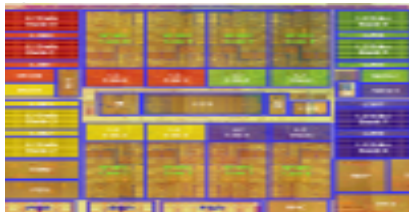
Intel Core i7
8 cores



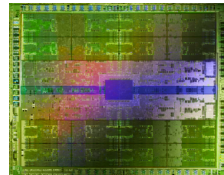
IBM Cell BE
8+1 cores



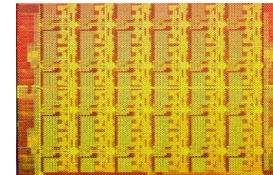
IBM POWER7
8 cores



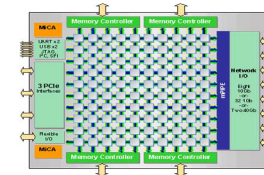
Sun Niagara II
8 cores



Nvidia Fermi
448 "cores"



Intel SCC
48 cores, networked

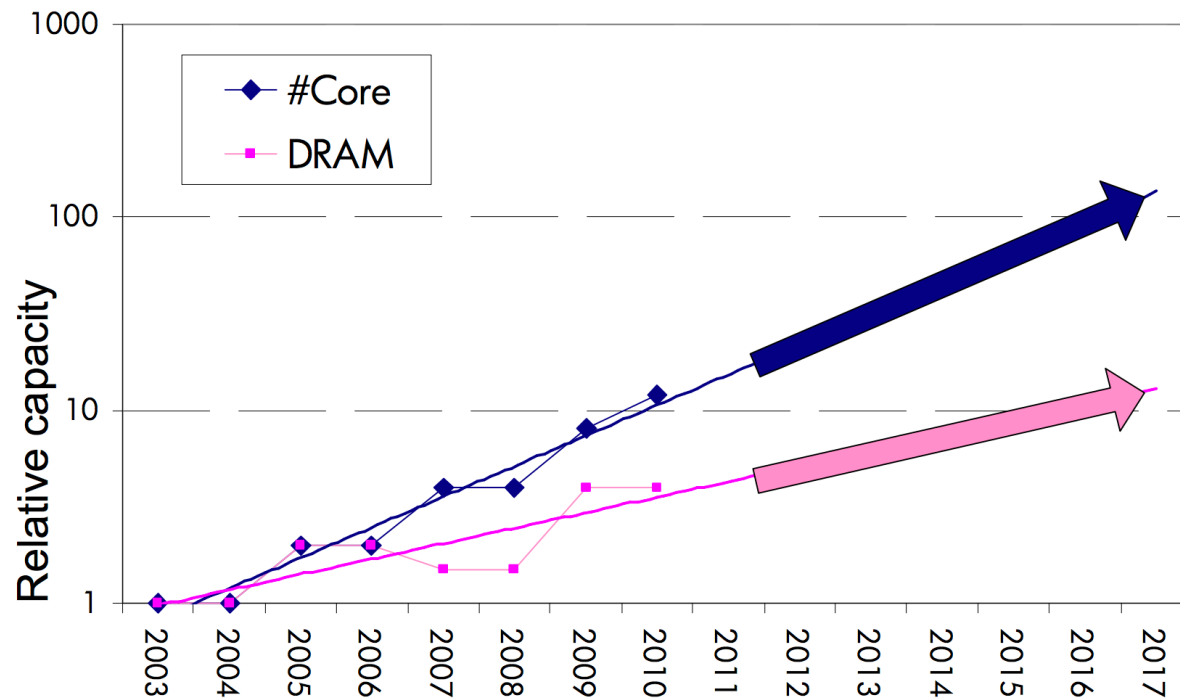


Tiler TILE Gx
100 cores, networked

CPU-DRAM Gap

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

Memory capacity per core expected to drop by 30% every two years



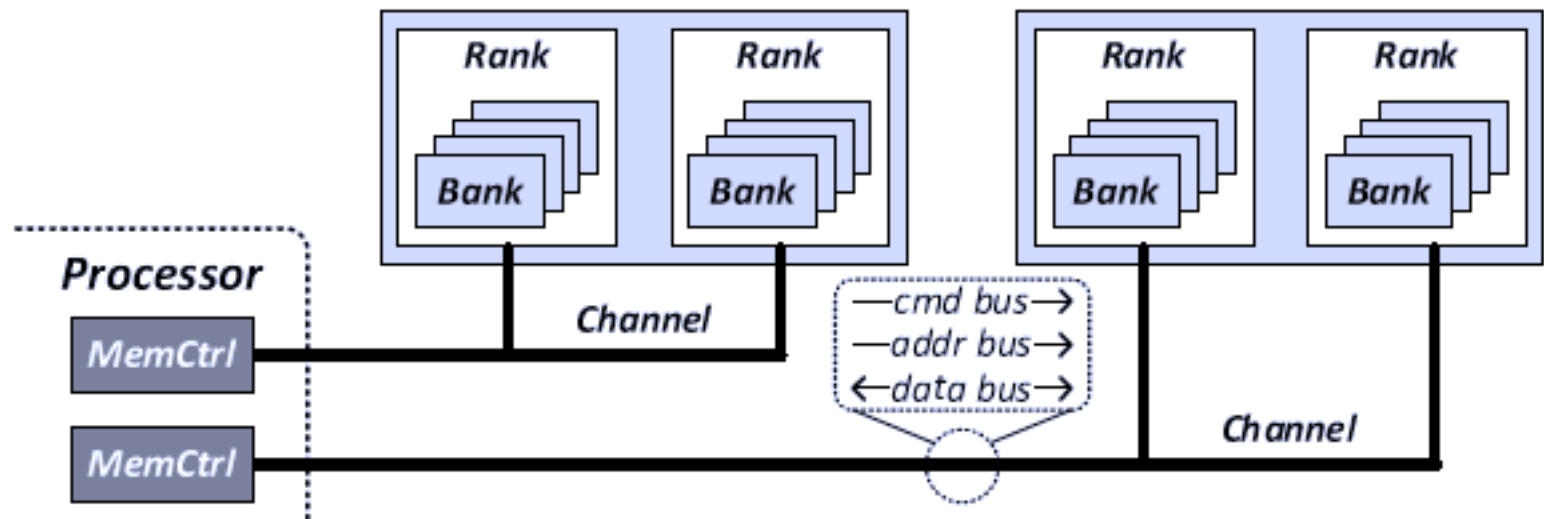
[Lim'09]

DRAM: Design Challenges

- Main memory energy/power is a key system design concern
 - Energy spent in off-chip memory hierarchy is about 40-50% [Lefurgy'03]
 - DRAM consumes power even when not used
 - periodic refresh
- DRAM technology scaling is ending
 - stops gaining higher capacity, lower cost, lower energy

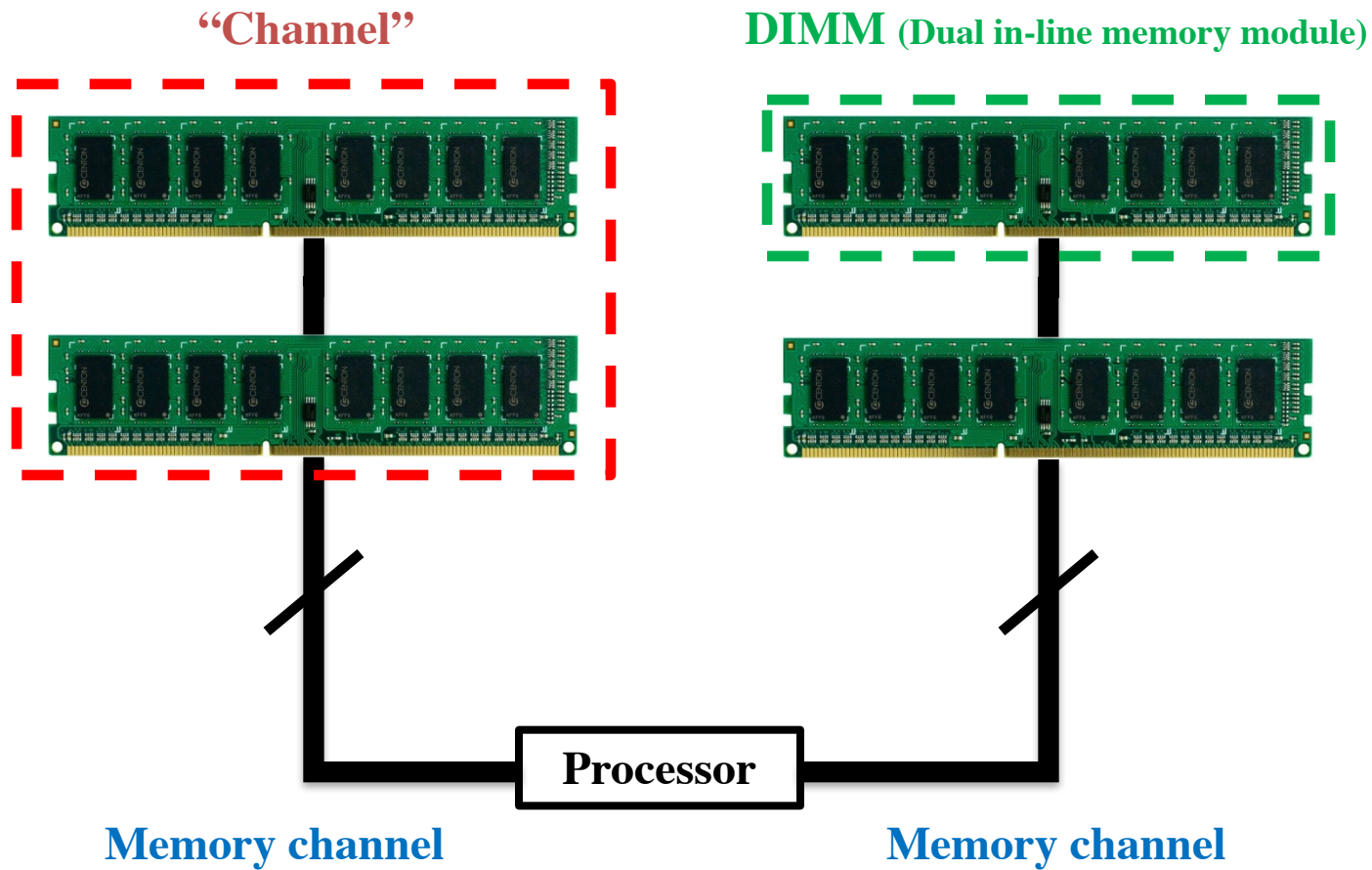
DRAM: Logical Organization

- Five DRAM coordinates
 - ▣ Channel, rank, bank, row, column



[Kim'12]

DRAM: Physical Organization



DIMM Structure

DIMM (Dual in-line memory module)



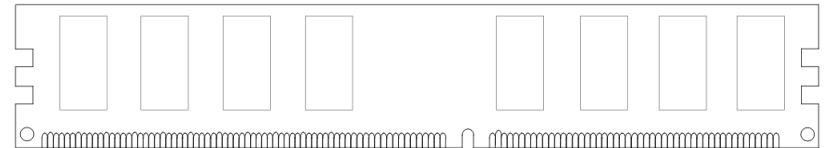
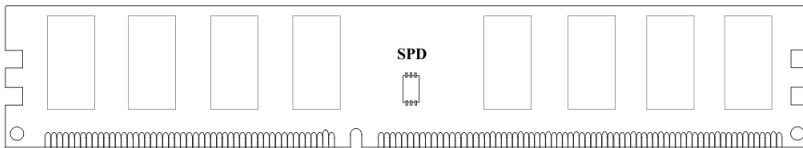
Side view

SIDE

4.00

Front of DIMM

Back of DIMM



DIMM Structure

DIMM (Dual in-line memory module)



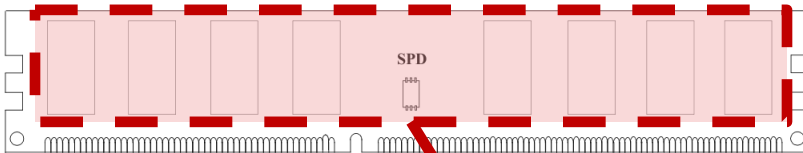
Side view

SIDE

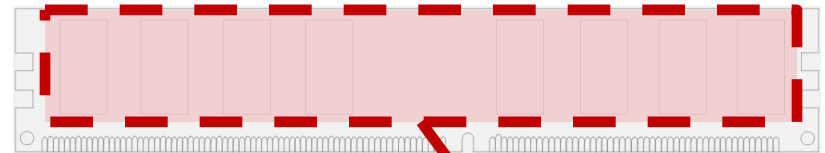
4.00

Front of DIMM

Back of DIMM



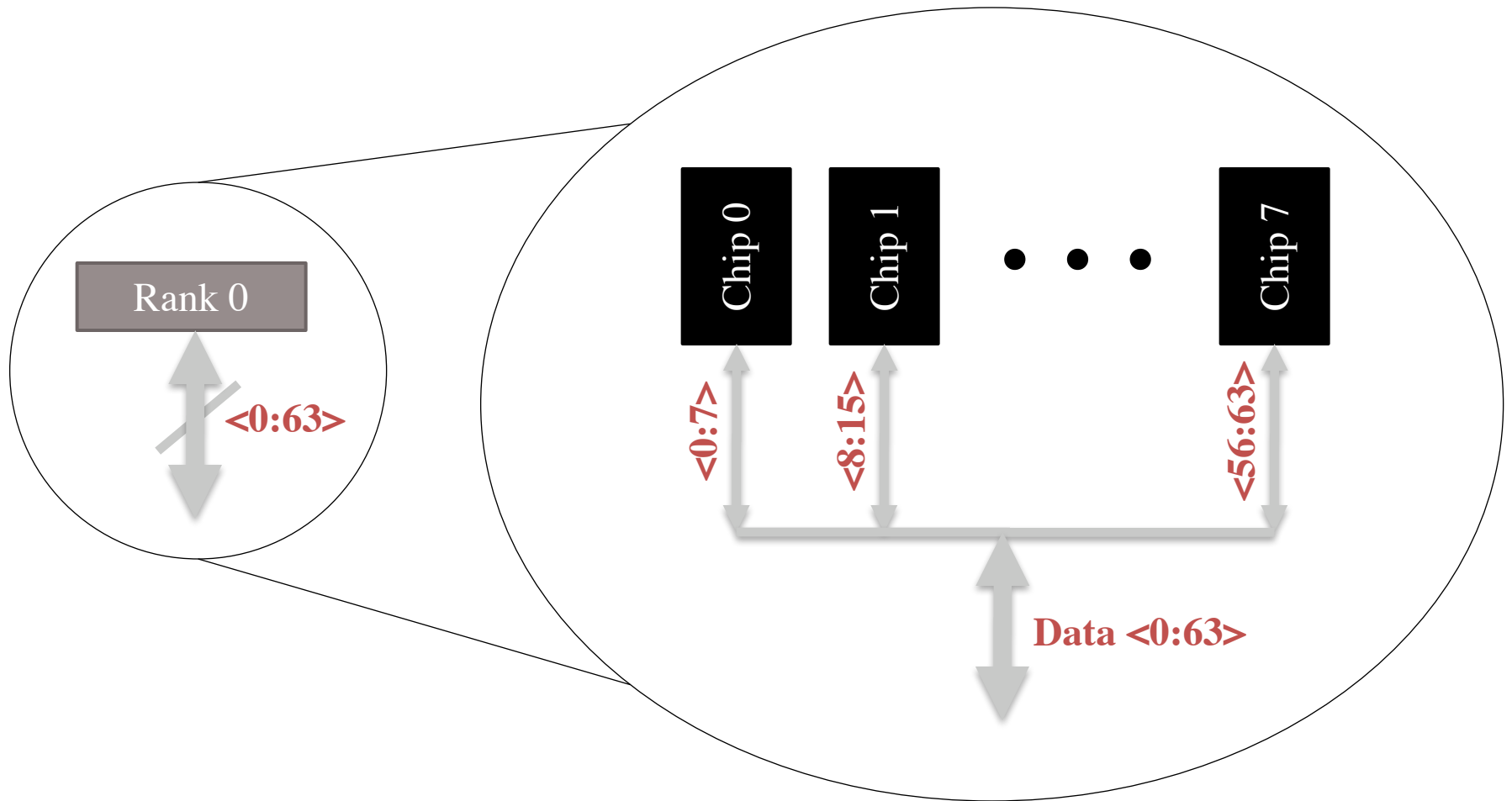
Rank 0: collection of 8 chips



Rank 1

Rank Organization

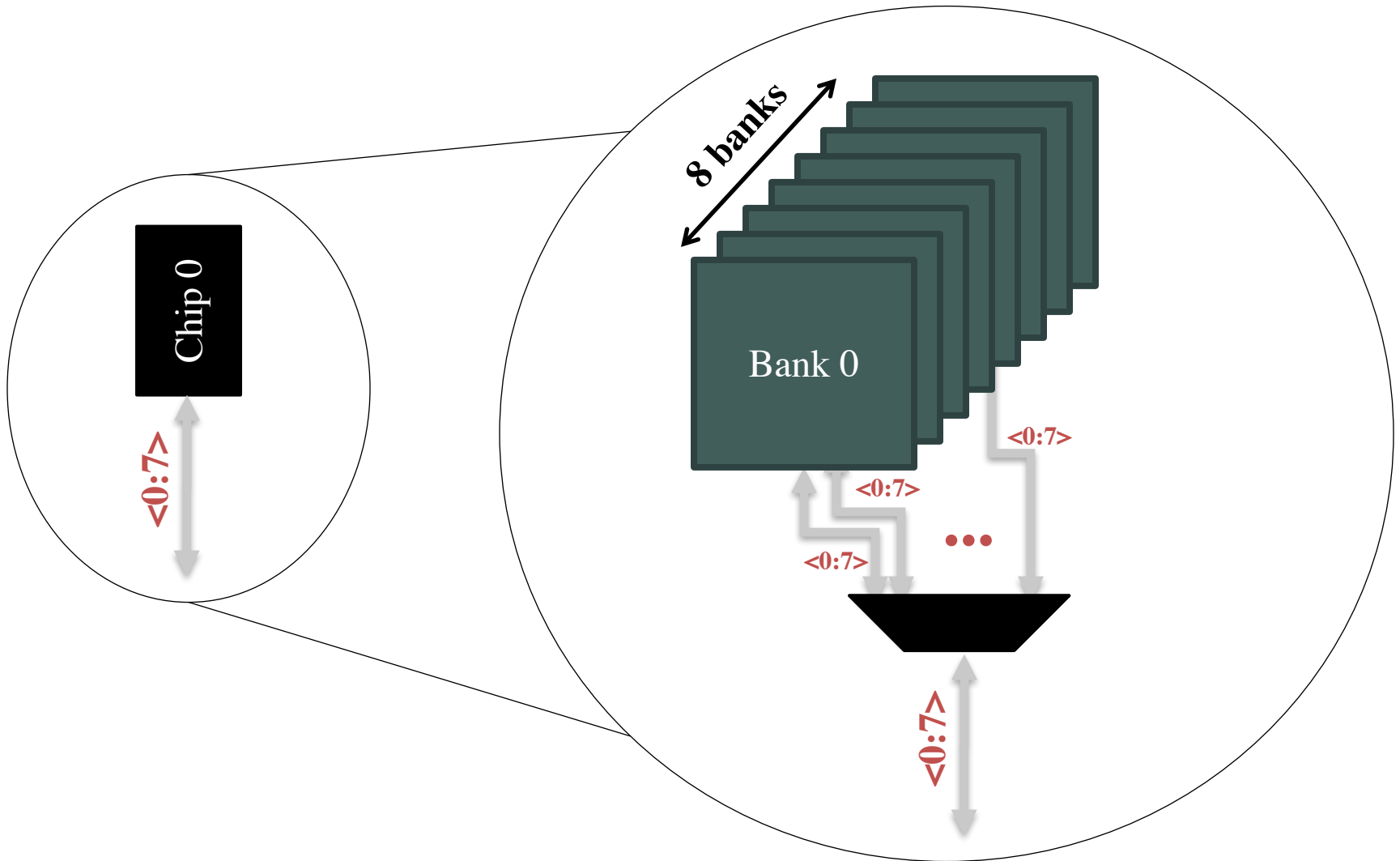
Rank Breakdown



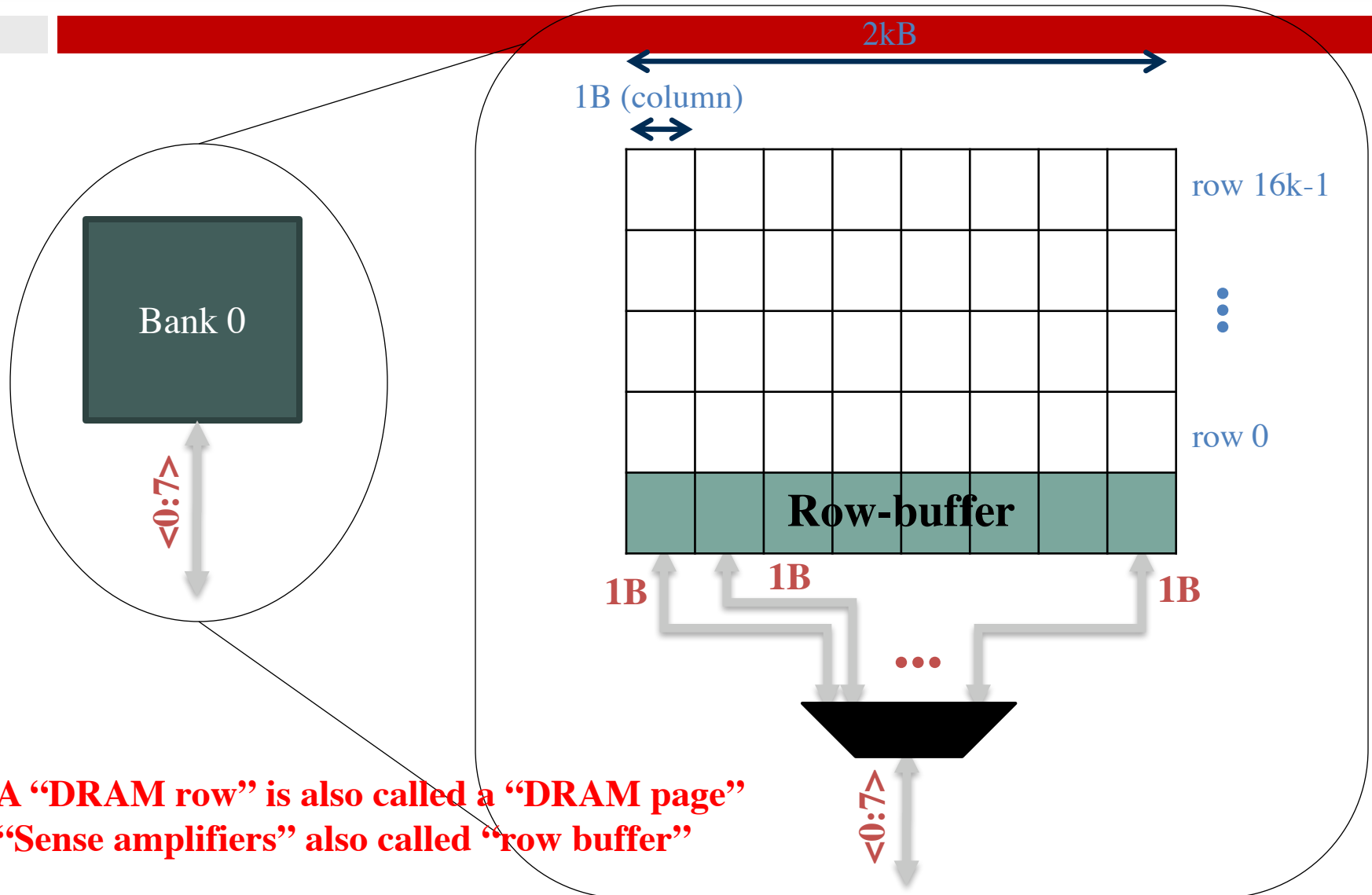
DRAM DIMM and Rank

- Multiple chips operated together to form a wide interface
 - ▣ All chips within a rank are controlled at the same time
 - ▣ Respond to a single command
 - ▣ Share address and command; different data bits
- A DRAM module consists of one or more ranks
 - ▣ e.g., DIMM (dual inline memory module)
 - ▣ If we have chips with 8-bit interface, to read 8 bytes in a single access, use 8 chips in a DIMM

Chip Structure



Bank Organization



DRAM Page Access

- 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 Page Access

Access Address:

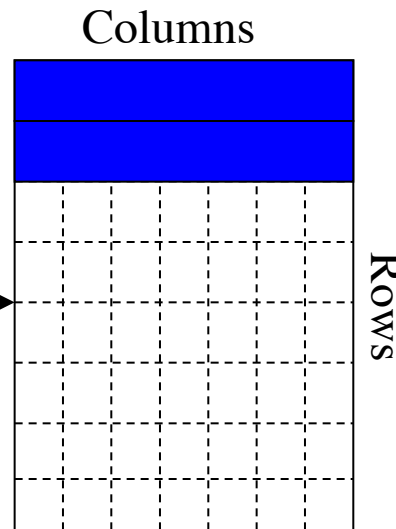
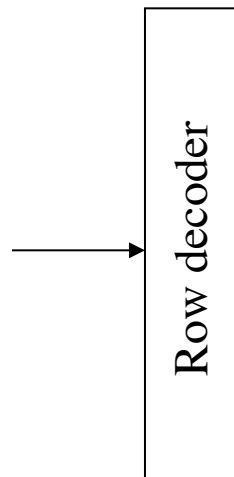
(Row 0, Column 0)

(Row 0, Column 1)

(Row 0, Column 85)

(Row 1, Column 0)

Row address 0



Row Buffer **CONFLICT !**

Column address 05

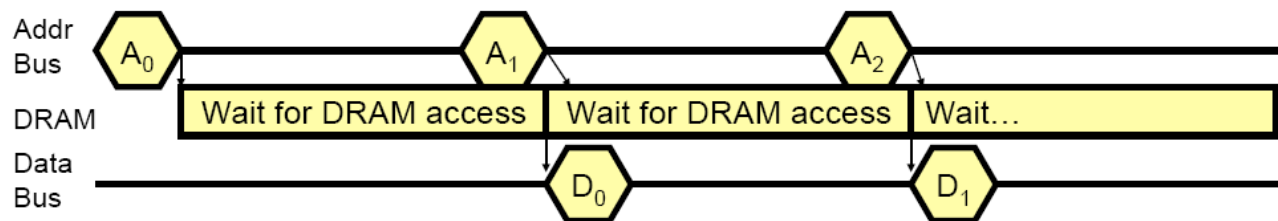


Data

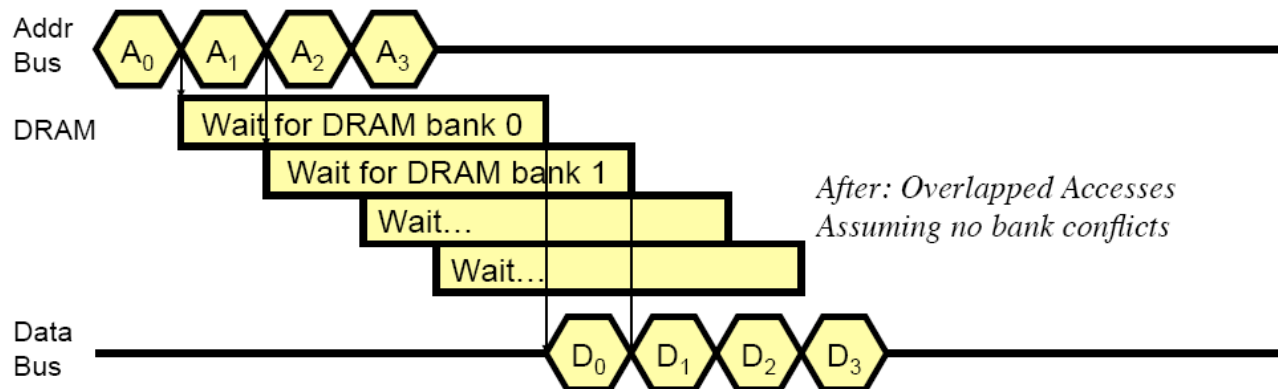
[ref: Mutlu]

DRAM: Parallel Access

- DRAM subsystem comprises multiple banks
 - ▣ Organized under independent channels and ranks



*Before: No Overlapping
Assuming accesses to different DRAM rows*



*After: Overlapped Accesses
Assuming no bank conflicts*

DRAM Controller

- 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

□ Ensuring DDRx timing constraints

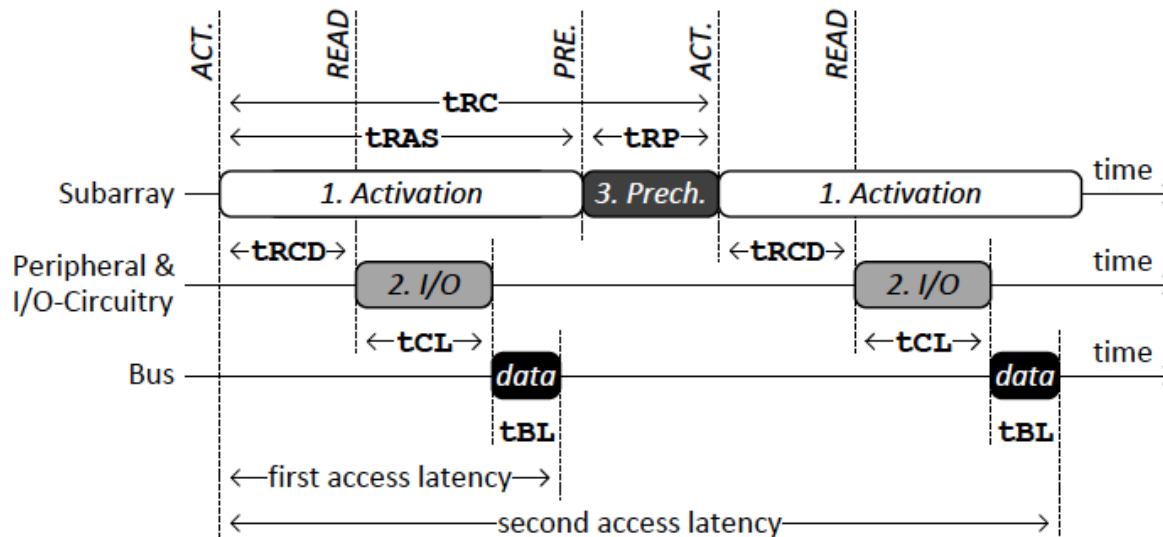


Figure 5. Three Phases of DRAM Access

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

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