

# Memory Systems

## Lec11 – Main Memory System

Chin-Fu Nien (粘徹夫)

(本節內容改自 Prof. Onur Mutlu, “Digital Design & Computer Architecture,” 22<sup>th</sup> Lecture, Spring 2011, “Introduction to Computer Architecture,” 21<sup>th</sup> Lecture, Spring 2015和“Computer Architecture,” 6<sup>th</sup>, 8<sup>th</sup>, 12<sup>th</sup> Lecture, Fall 2023 講義)

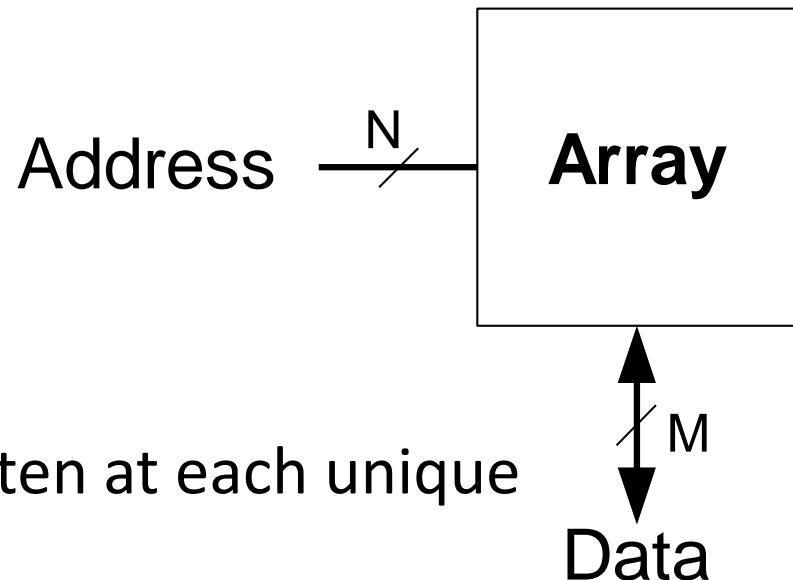
# Module 3: Advanced Architectural Topics

# Quick Overview of Memory Arrays

(本節內容改自 Prof. Onur Mutlu, “Digital Design & Computer Architecture,” 22<sup>th</sup> Lecture, Spring 2021課程講義)

# Array Organization of Memories

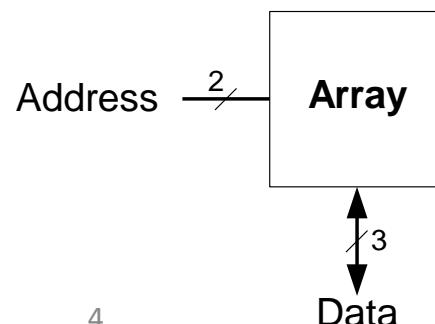
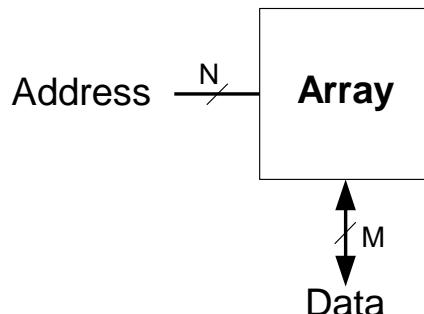
- Goal: Efficiently store large amounts of data
  - A memory array (stores data)
  - Address selection logic (selects one row of the array)
  - Readout circuitry (reads data out)



- An M-bit value can be read or written at each unique N-bit address
  - All values can be accessed, but only M-bits at a time
  - Access restriction allows more compact organization

# Memory Arrays

- Two-dimensional array of bit cells
  - Each bit cell stores one bit
- An array with N address bits and M data bits:
  - $2^N$  rows and M columns
  - Depth: number of rows (number of words)
  - Width: number of columns (size of word)
  - Array size: depth  $\times$  width =  $2^N \times M$



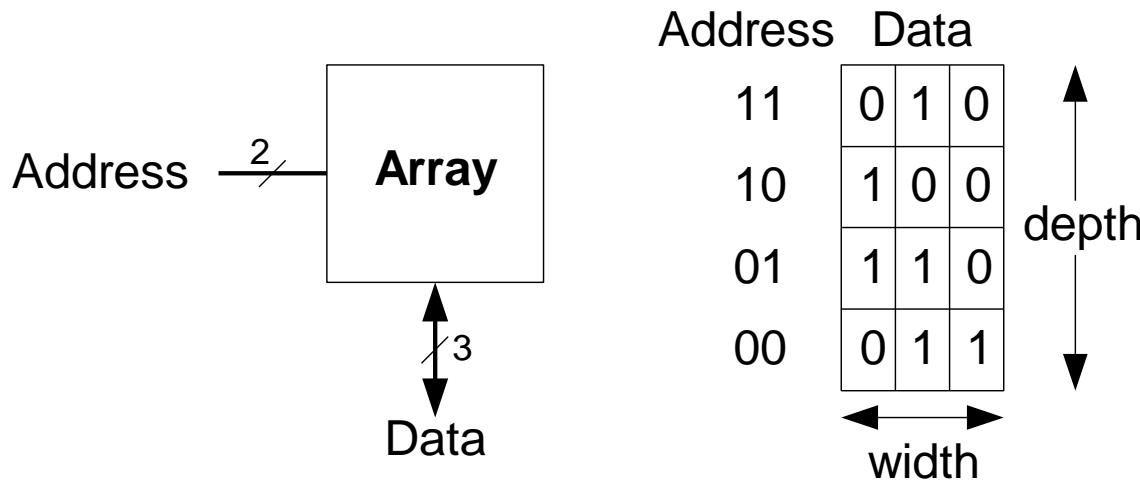
	Address	Data
11	0 1 0	
10	1 0 0	
01	1 1 0	
00	0 1 1	

depth

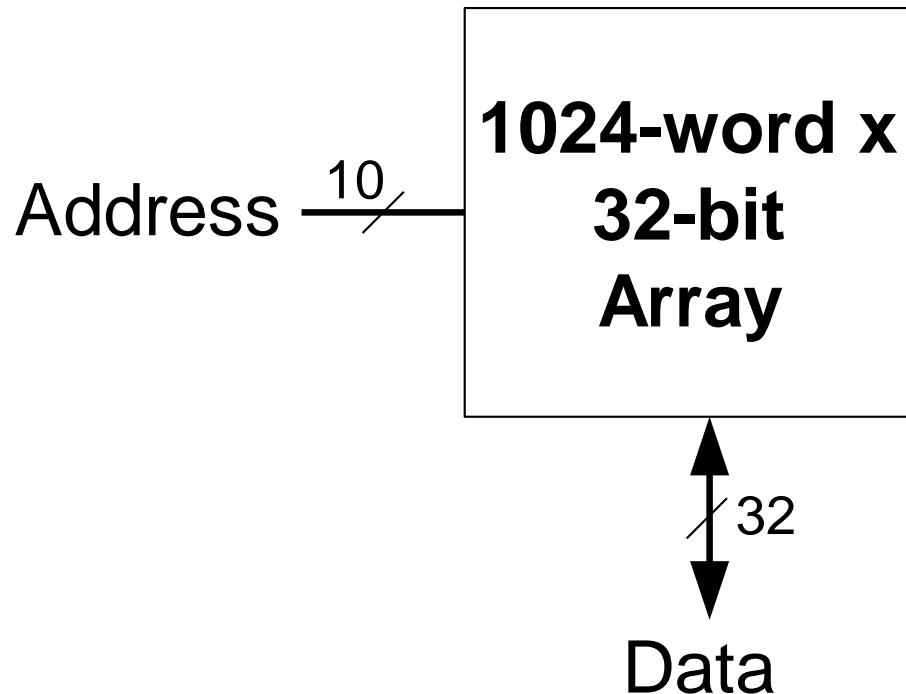
width

# Memory Array Example

- $2^2 \times 3$ -bit array
- Number of words: 4
- Word size: 3-bits
- For example, the 3-bit word stored at address 10 is 100

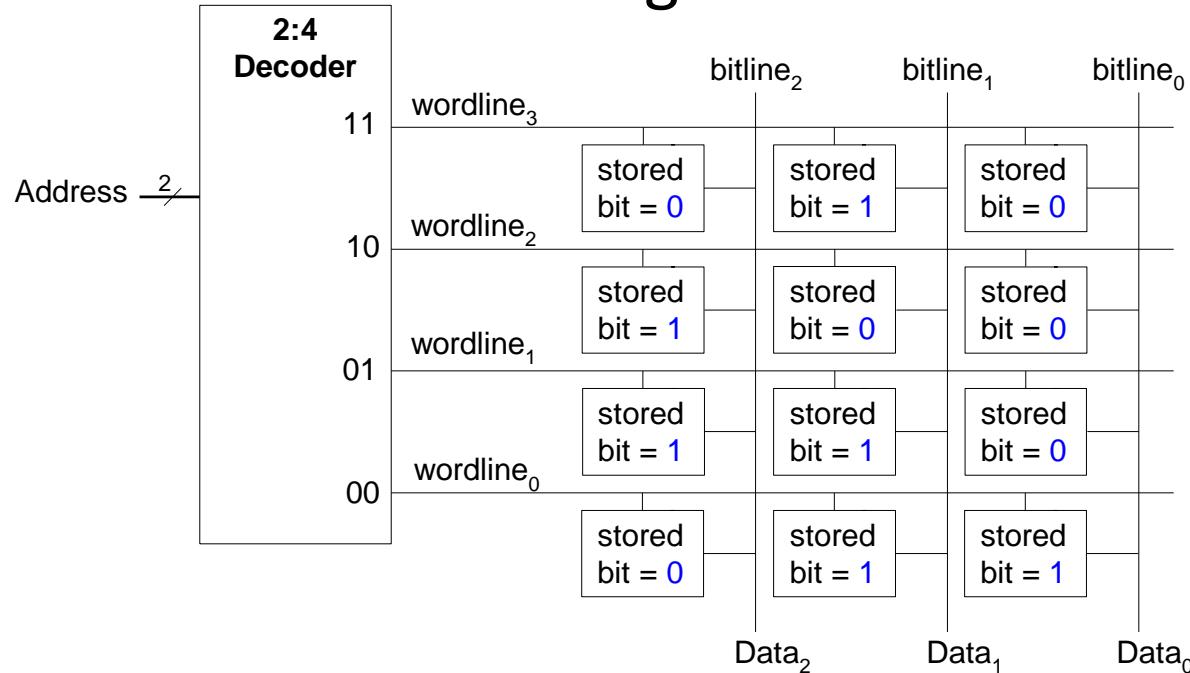


# Larger and Wider Memory Array Example



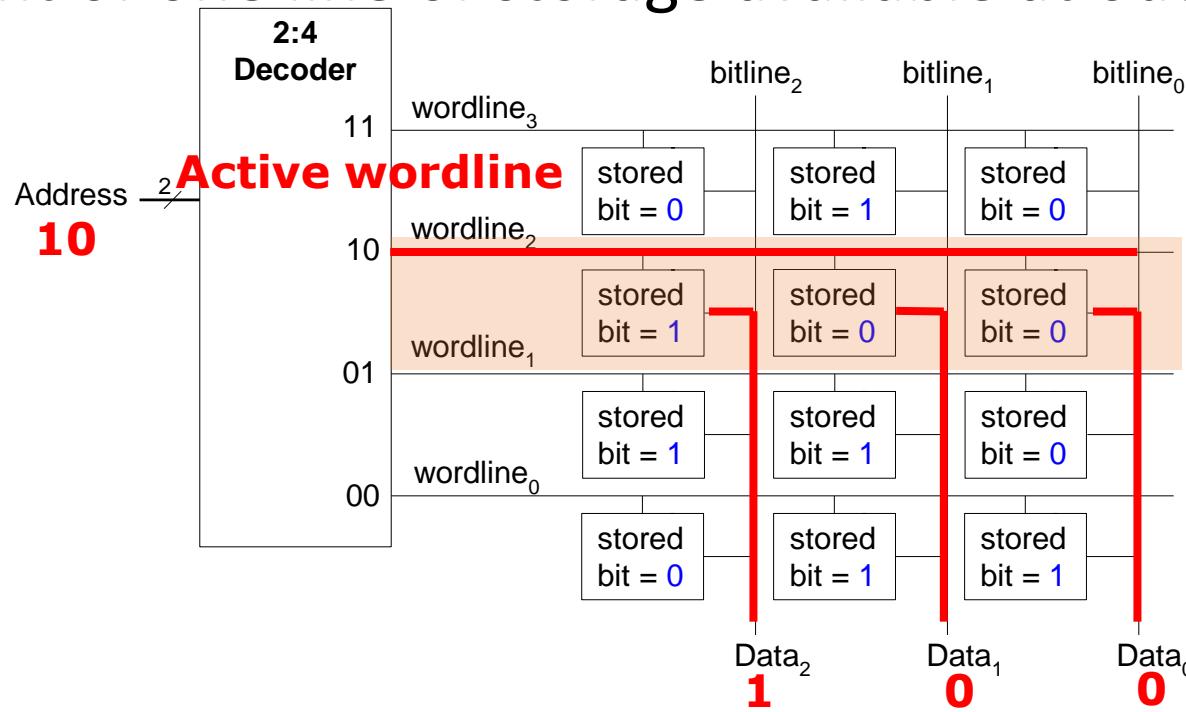
# Memory Array Organization (I)

- Storage nodes in one column connected to one bitline
- Address decoder activates only ONE wordline
- Content of one line of storage available at output

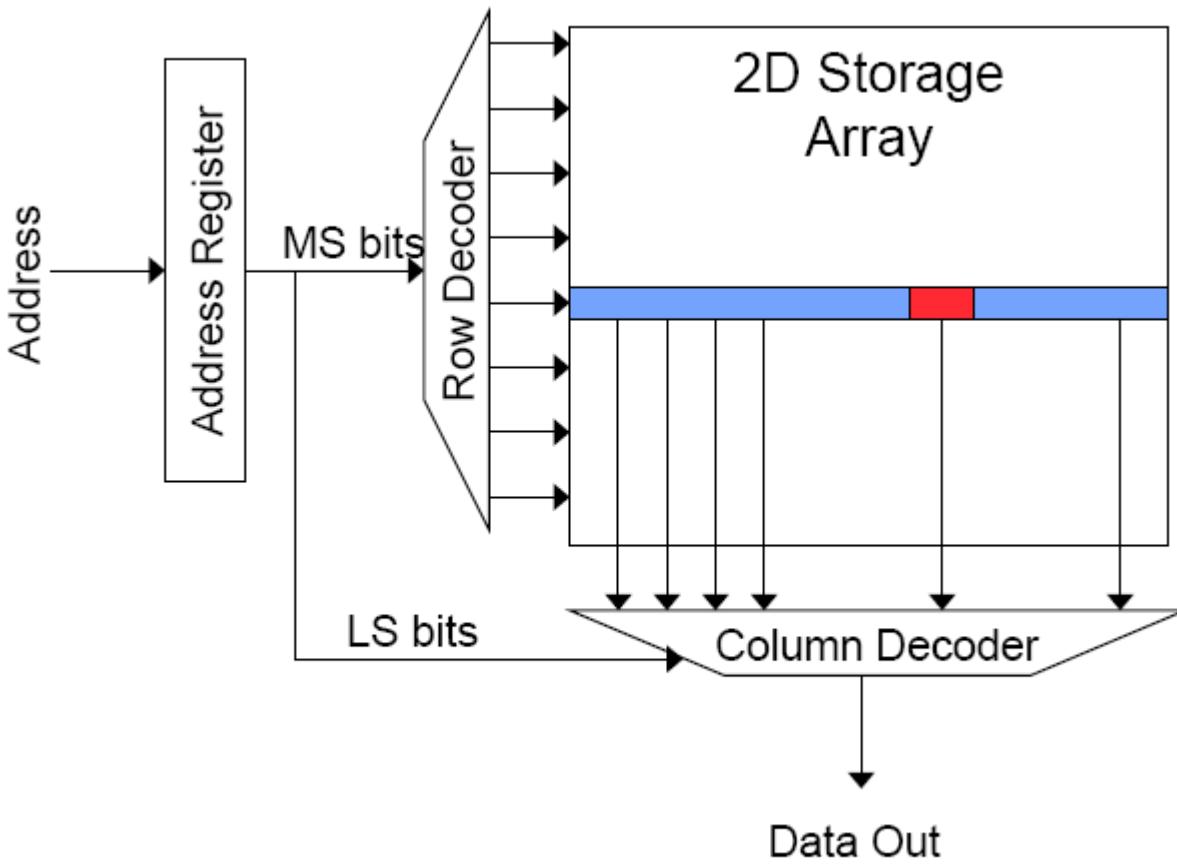


# Memory Array Organization (II)

- Storage nodes in one column connected to one bitline
- Address decoder activates only ONE wordline
- Content of one line of storage available at output

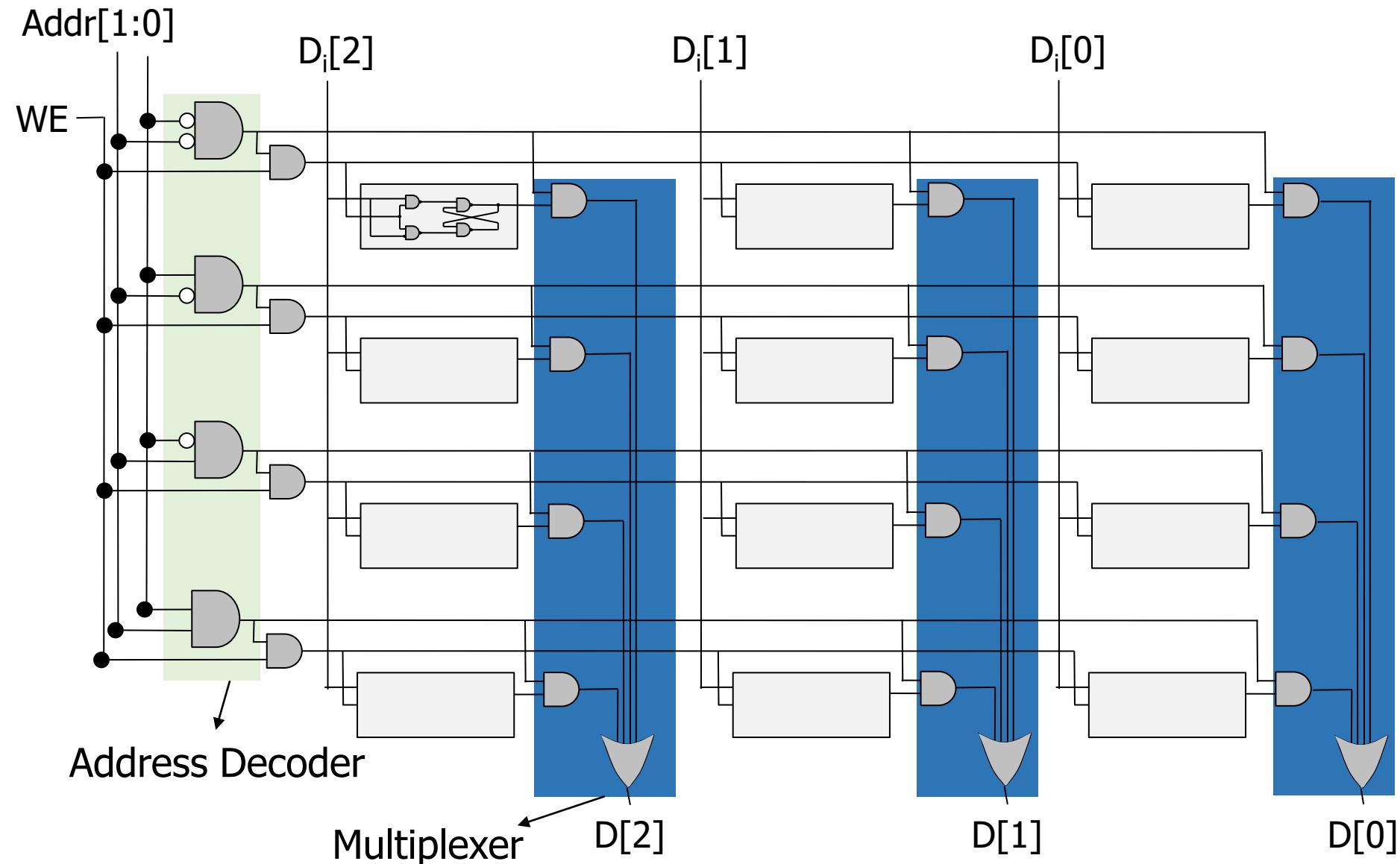


# Memory Bank Organization and Operation

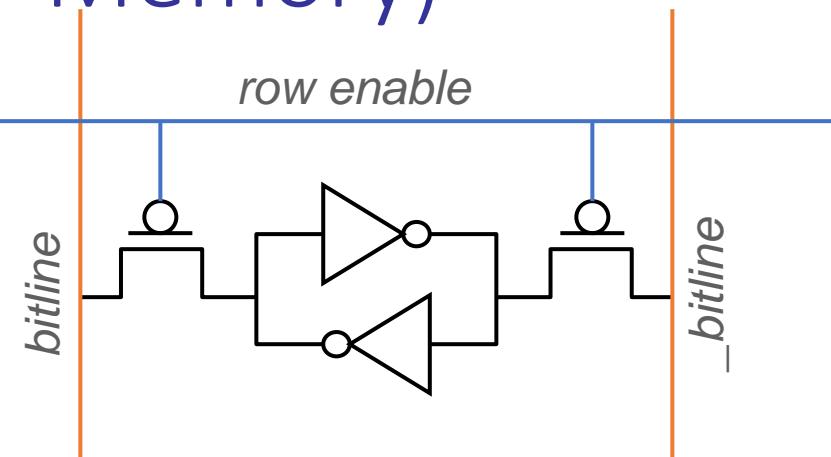


- Read access sequence:
  1. Decode row address & drive word-lines
  2. Selected bits drive bit-lines
    - Entire row read
  3. Amplify row data
  4. Decode column address & select subset of row
    - Send to output
  5. Precharge bit-lines
    - For next access

# A Bigger Memory Array (4 locations X 3 bits)

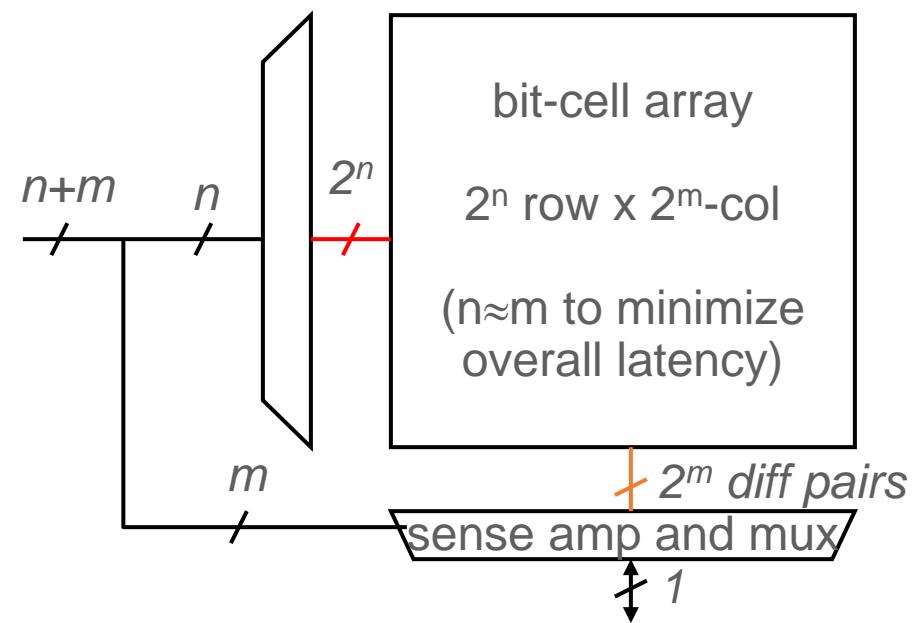


# SRAM (Static Random Access Memory)



## Read Sequence

1. address decode
2. drive row select
3. selected bit-cells drive bitlines  
(entire row is read together)
4. differential sensing and column select  
(data is ready)
5. precharge all bitlines  
(for next read or write)

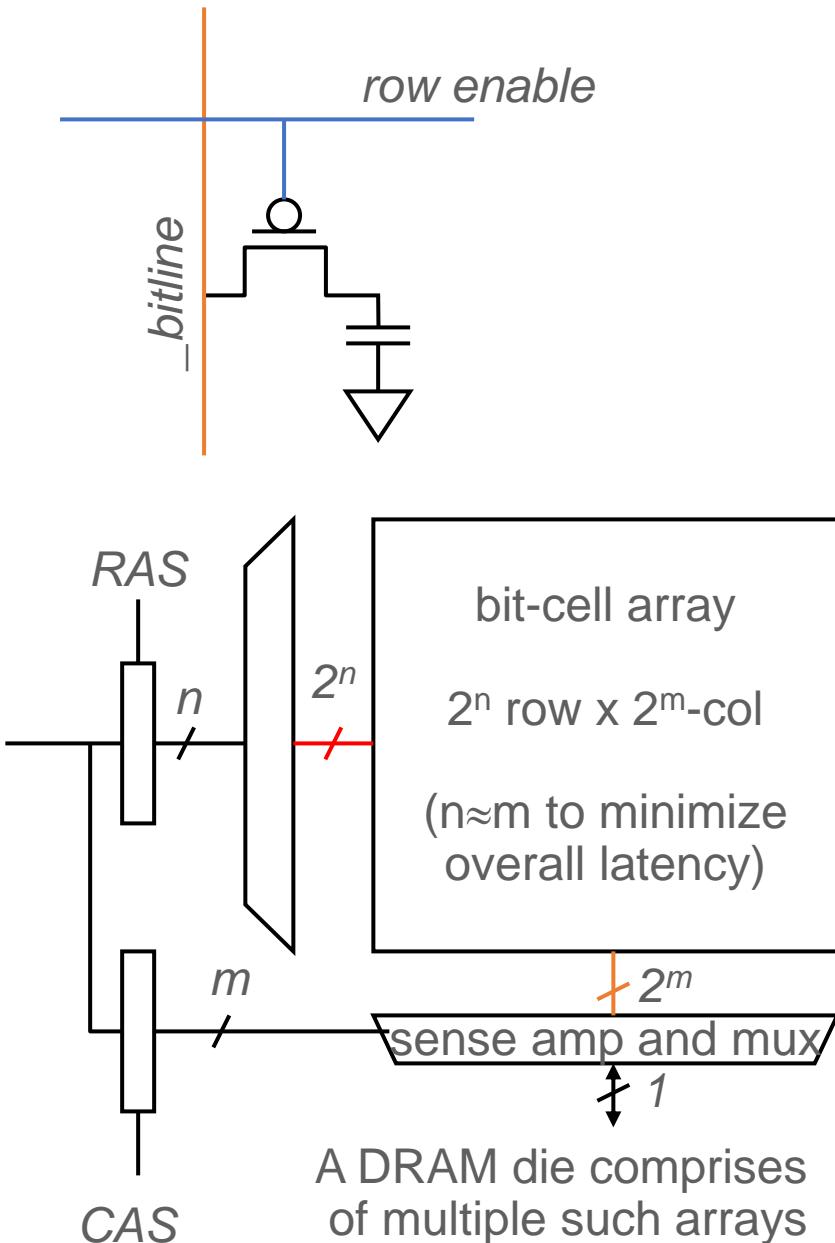


Access latency dominated by steps 2 and 3

Cycling time dominated by steps 2, 3 and 5

- step 2 proportional to  $2^m$
- step 3 and 5 proportional to  $2^n$

# DRAM (Dynamic Random Access Memory)



Bit stored as charge on node capacitor (non-restorative)

- bit cell loses charge when read
- bit cell loses charge over time

Read Sequence

- 1~3 same as SRAM
4. a “flip-flopping” sense amp amplifies and regenerates the bitline, data bit is mux’ ed out
5. precharge all bitlines

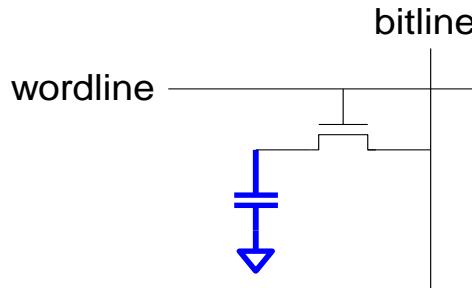
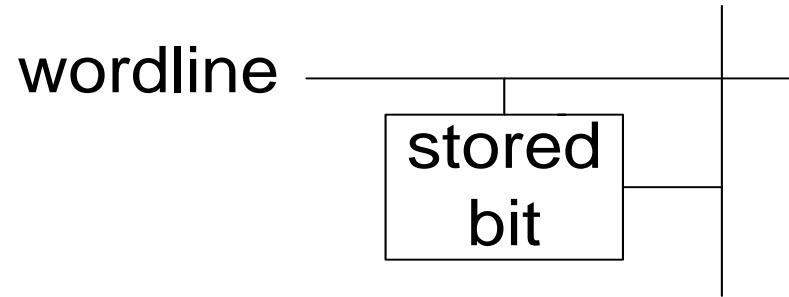
Destructive reads

Charge loss over time

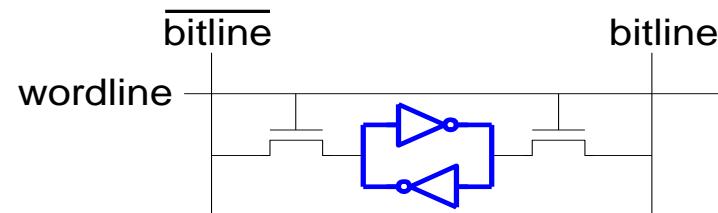
**Refresh:** A DRAM controller must periodically read each row within the allowed refresh time (10s of ms) such that charge is restored

# How is Access Controlled?

- Access transistors (that are configured as switches) connect the bit storage to the bitline
- Access controlled by the wordline



**DRAM**



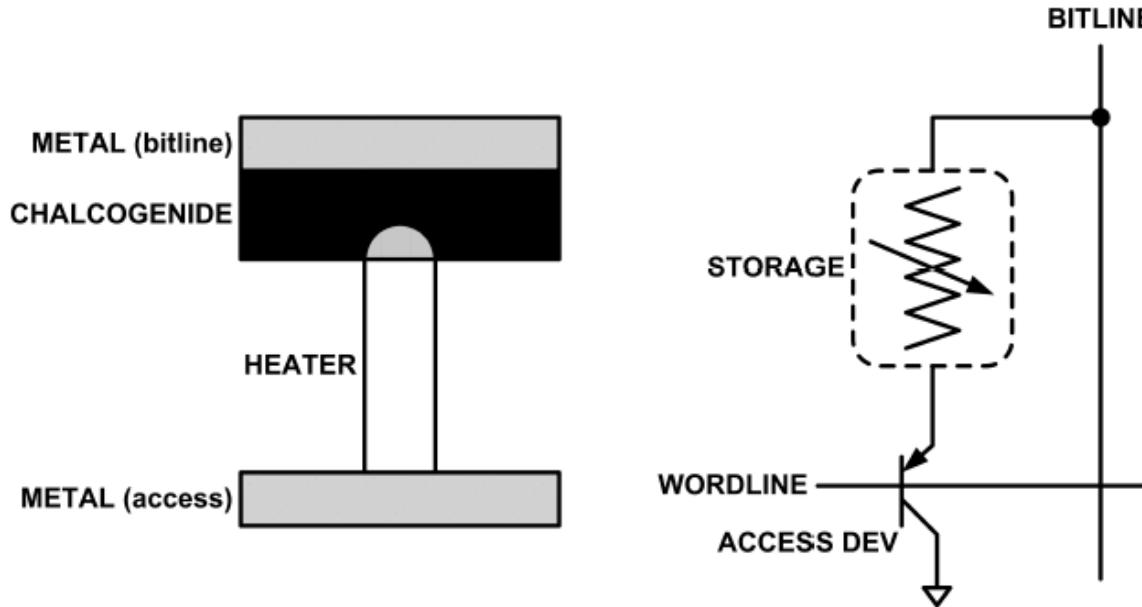
**SRAM**

# DRAM vs. SRAM

- DRAM
  - Slower access (capacitor)
  - Higher density (1T 1C cell)
  - Lower cost
  - Requires refresh (power, performance, circuitry)
  - Manufacturing requires putting capacitor and logic together
- SRAM
  - Faster access (no capacitor)
  - Lower density (6T cell)
  - Higher cost
  - No need for refresh
  - Manufacturing compatible with logic process (no capacitor)

# 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



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.

# Intel Optane Persistent Memory (2019)

- Non-volatile main memory
- Based on 3D-XPoint Technology



# DRAM vs. PCM

- DRAM

- Faster access (capacitor)
- Lower density (capacitor less scalable) → higher cost
- Requires refresh (power, performance, circuitry)
- Manufacturing requires putting capacitor and logic together
- Volatile (loses data at loss of power)
- No endurance problems
- Lower access energy

- PCM

- Slower access (no capacitor)
- Higher density (phase change material more scalable) → lower cost
- No need for refresh
- Manufacturing requires less conventional processes – less mature
- Non-volatile (does not lose data at loss of power)
- Endurance problems (a cell cannot be used after N writes)
- Higher access energy

# Building Larger Memories

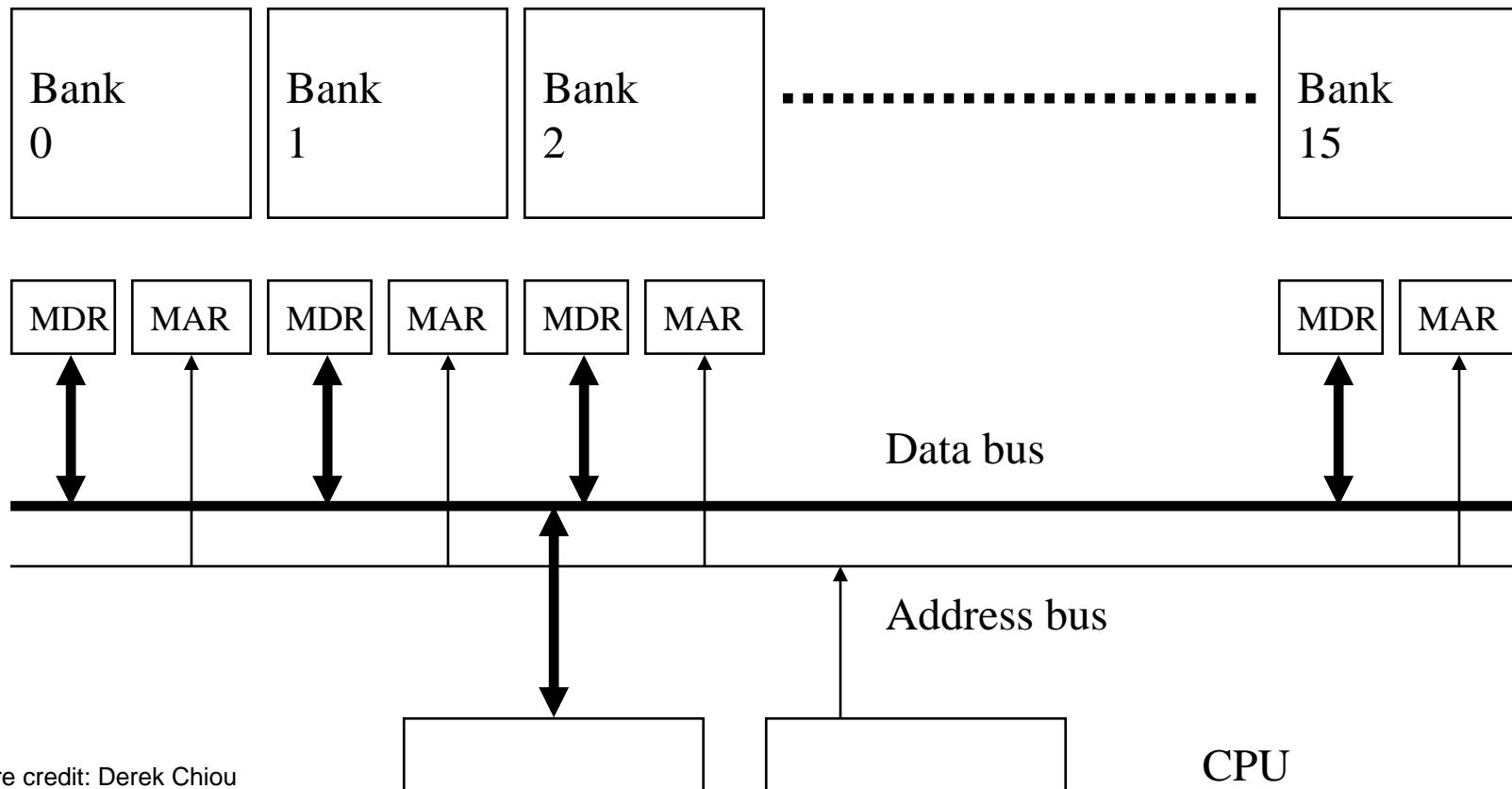
- Requires larger memory arrays
- Large → slow
- How do we make the memory large without making it too slow?
- Idea: Divide the memory into smaller arrays and interconnect the arrays to input/output buses
  - Large memories are hierarchical array structures
  - DRAM: Channel → Rank → Bank → Subarrays → Mats

# General Principle: Interleaving (Banking)

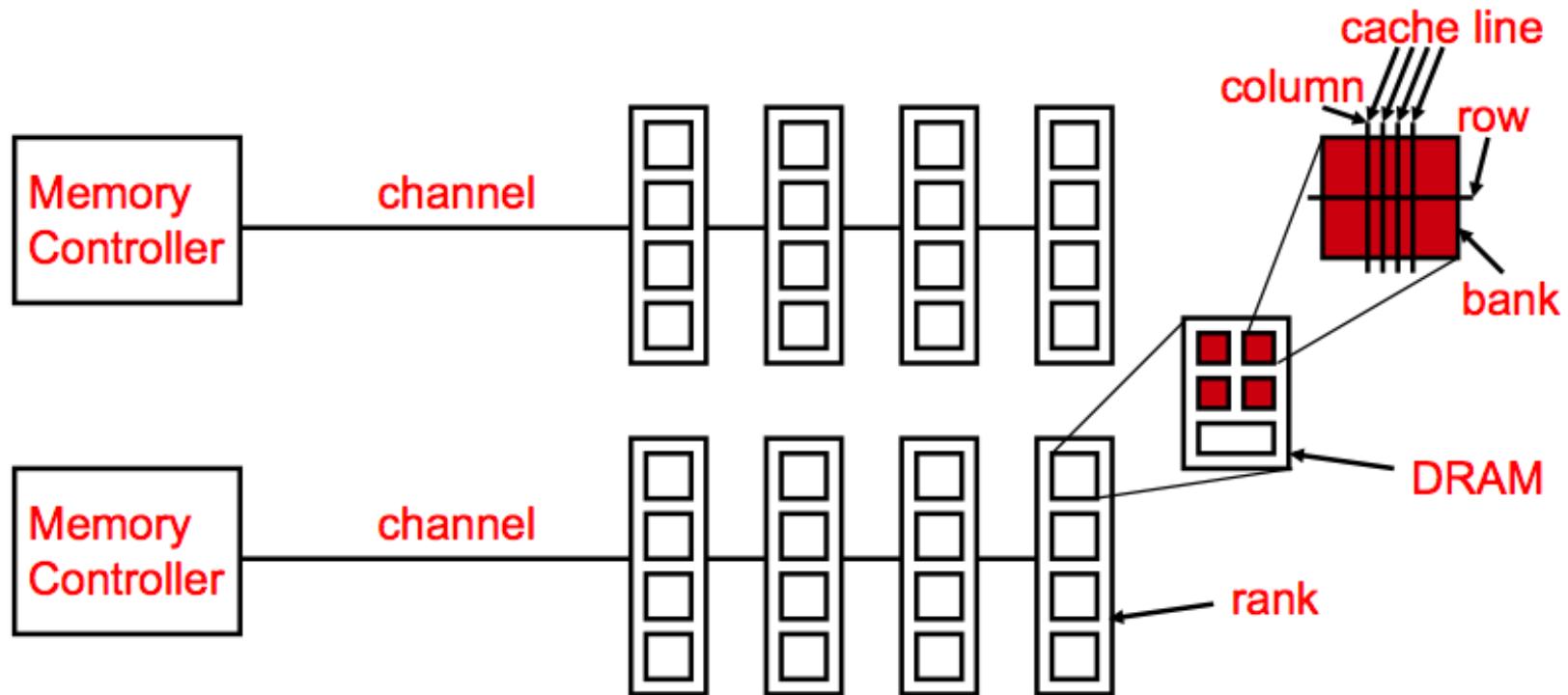
- **Interleaving (banking)**
  - **Problem:** a single monolithic large memory array takes long to access and does not enable multiple accesses in parallel
  - **Goal:** Reduce the latency of memory array access and enable multiple accesses in parallel
  - **Idea:** Divide a large array into multiple 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
  - **A Key Issue:** How do you map data to different banks? (i.e., how do you interleave data across banks?)

# Recall: Memory Banking

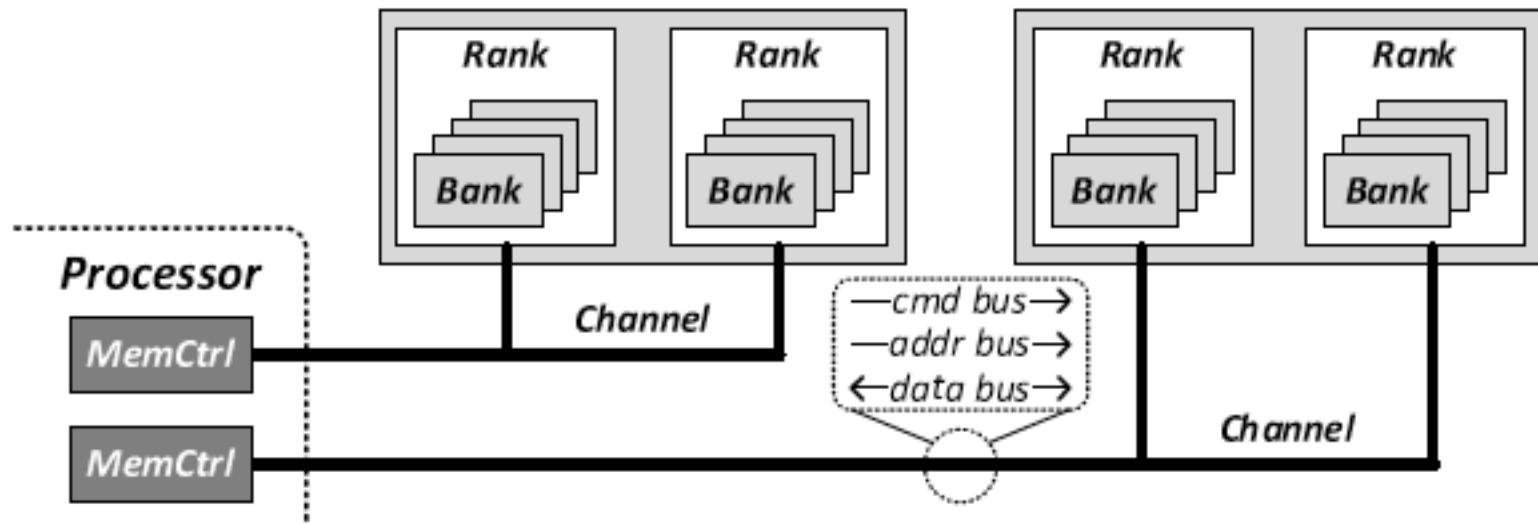
- Memory is divided into **banks** that can be accessed independently; banks share address and data buses (to minimize pin cost)
- Can start and complete one bank access per cycle
- **Can sustain N concurrent accesses if all N go to different banks**



# Generalized Memory Structure



# Generalized Memory Structure



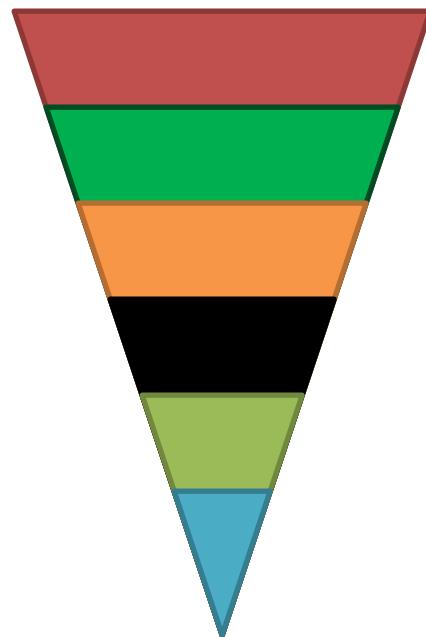
Kim+, “A Case for Exploiting Subarray-Level Parallelism in DRAM,” ISCA 2012.  
Lee+, “Decoupled Direct Memory Access,” PACT 2015.

# The DRAM System

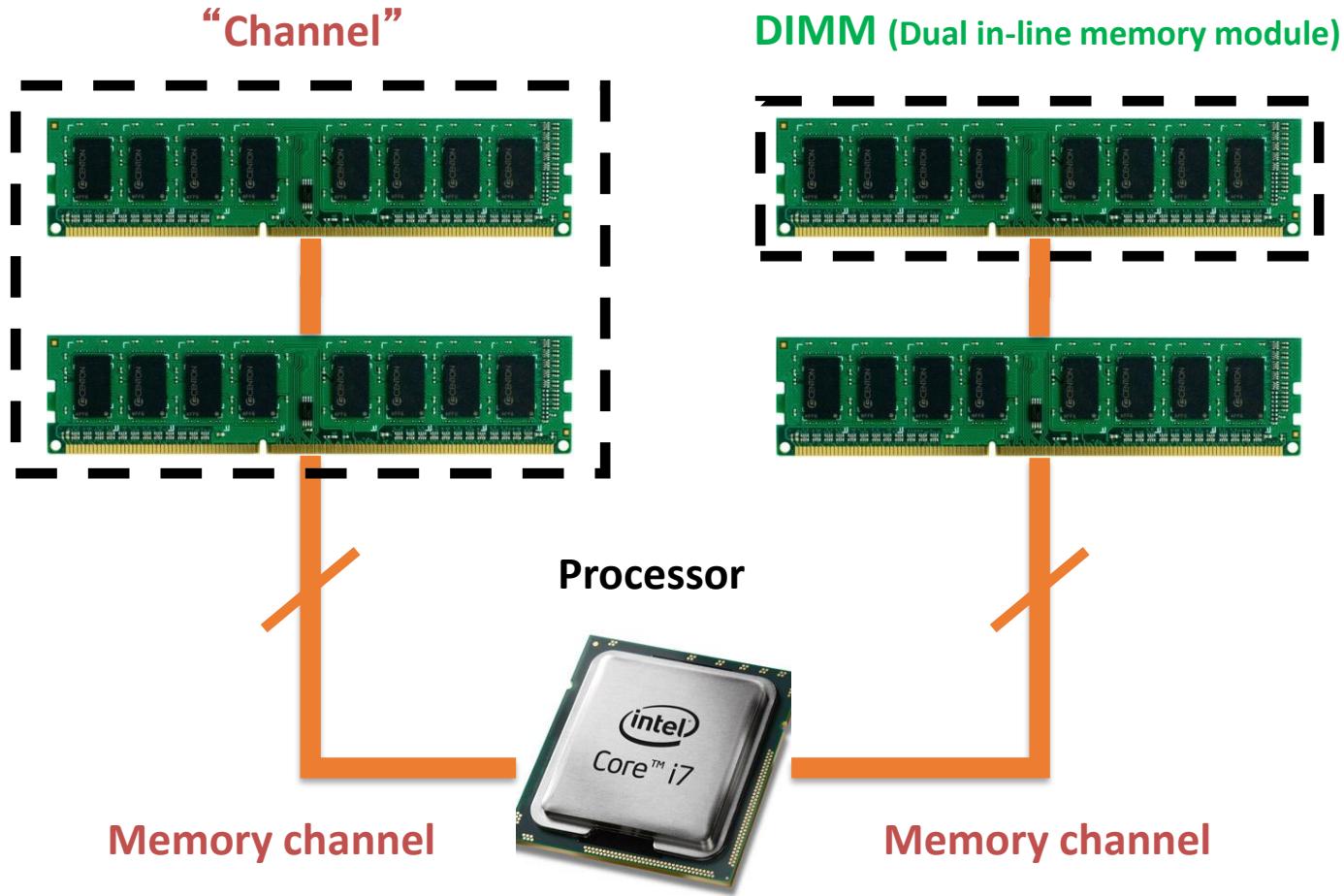
(本節內容改自 Prof. Onur Mutlu, “Digital Design & Computer Architecture,” 22<sup>th</sup> Lecture, Spring 2021與“Introduction to Computer Architecture,” 21<sup>th</sup> Lecture, Spring 2015課程講義)

# DRAM Subsystem Organization

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

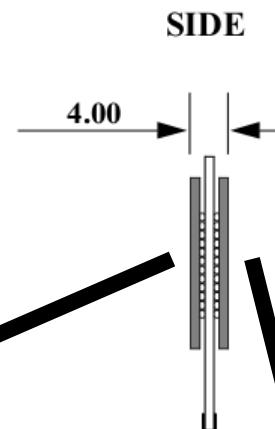


# The DRAM Subsystem

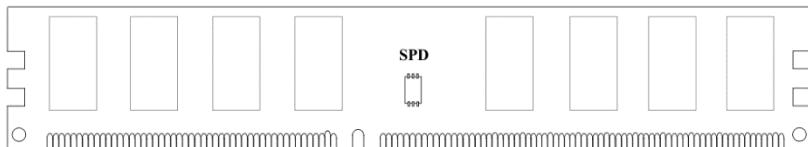


# Breaking down a DIMM (module)

DIMM (Dual in-line memory module)



Front of DIMM



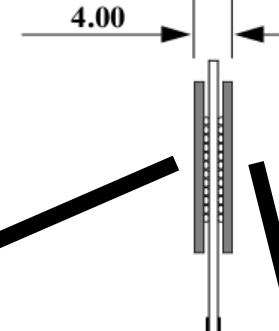
Back of DIMM

# Breaking down a DIMM (module)

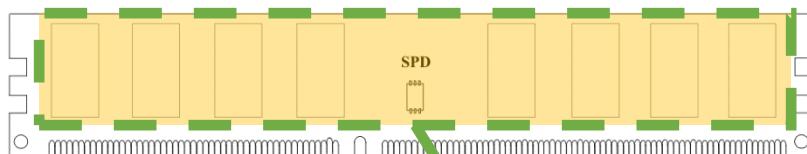
DIMM (Dual in-line memory module)



SIDE

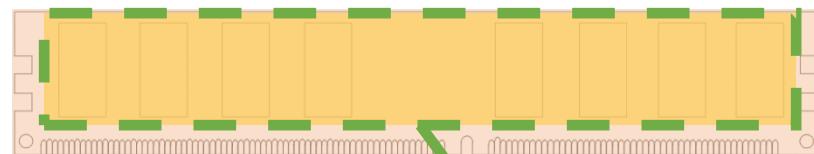


Front of DIMM



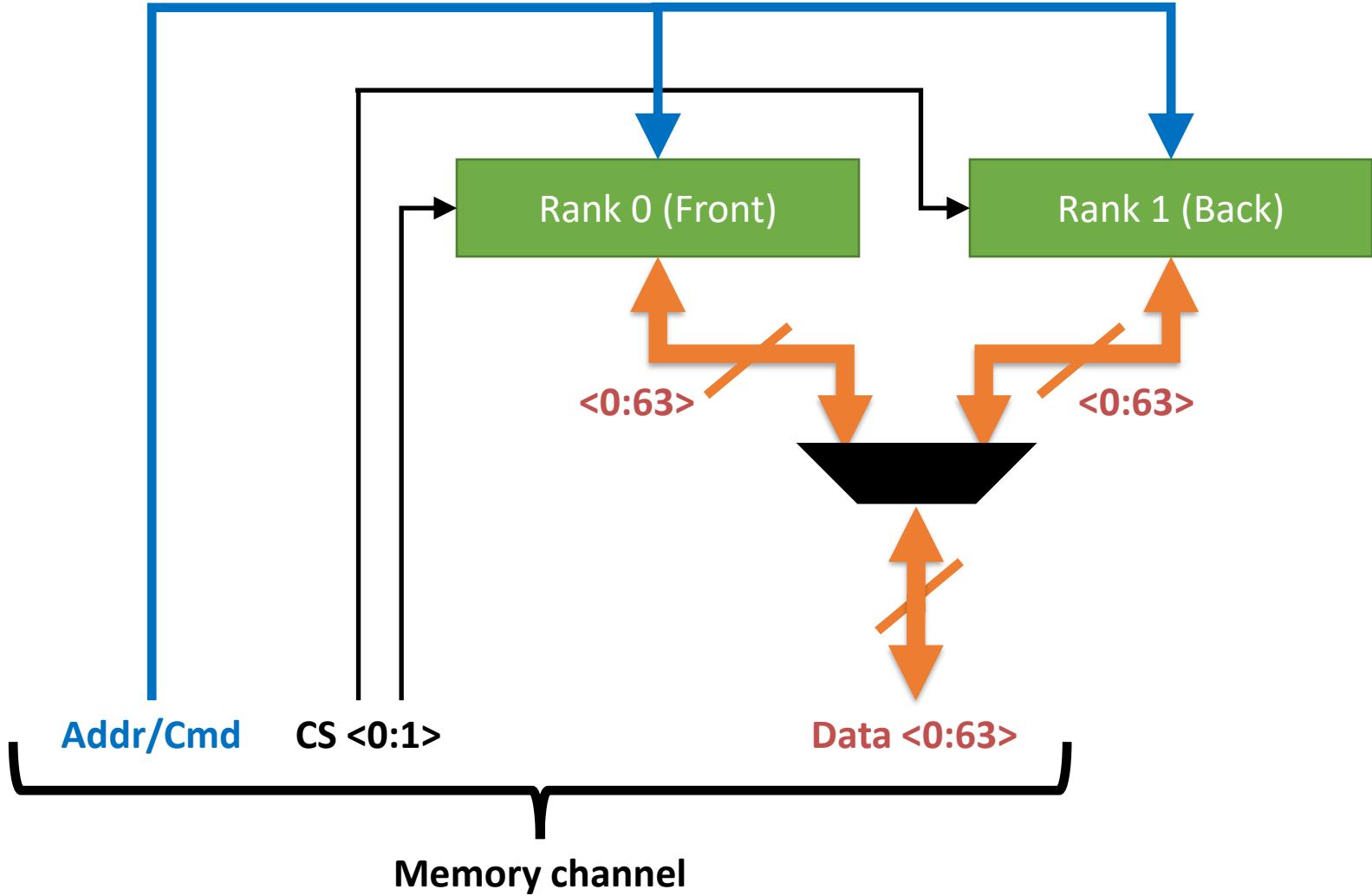
Rank 0: collection of 8 chips

Back of DIMM

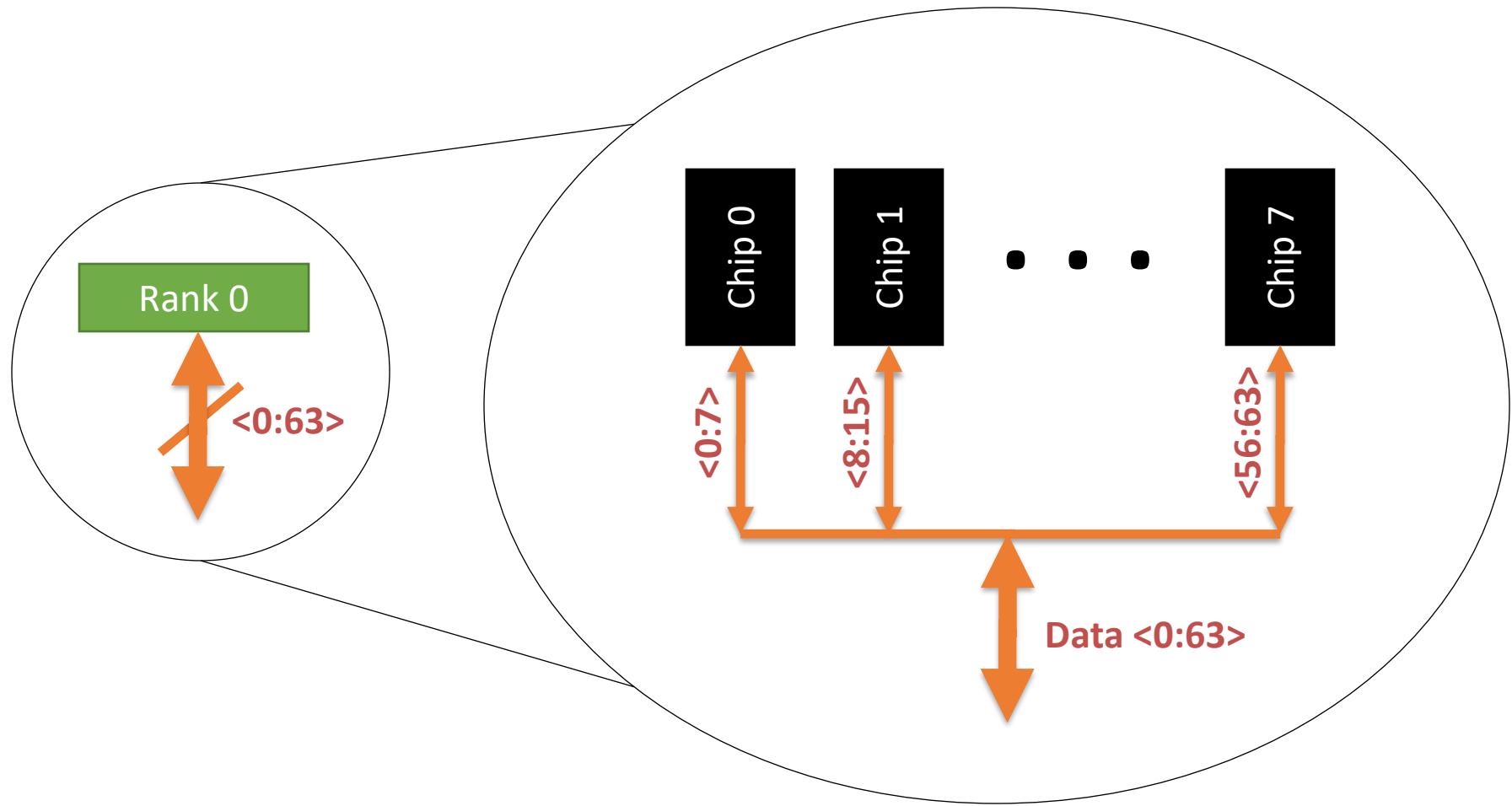


Rank 1

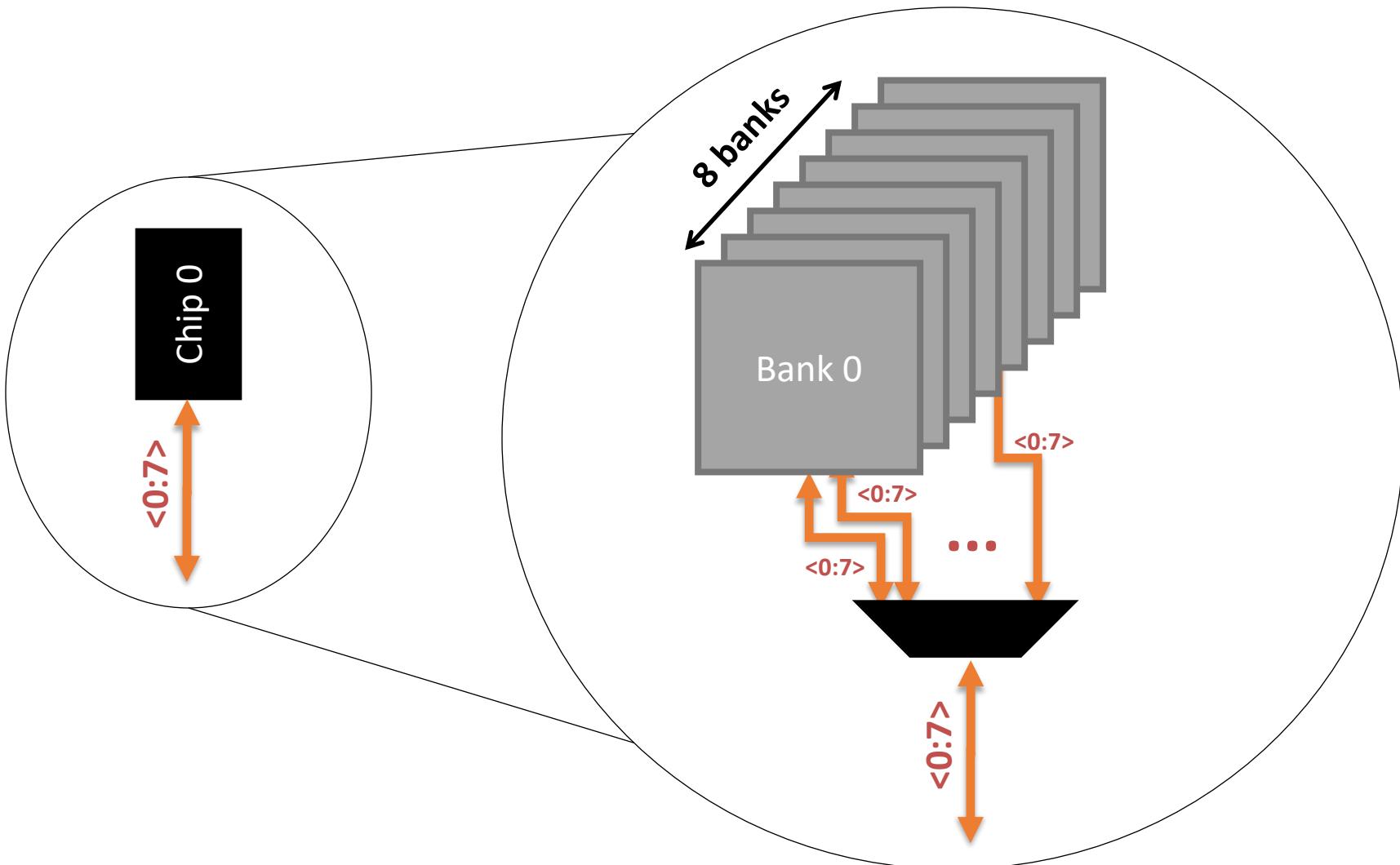
# Rank



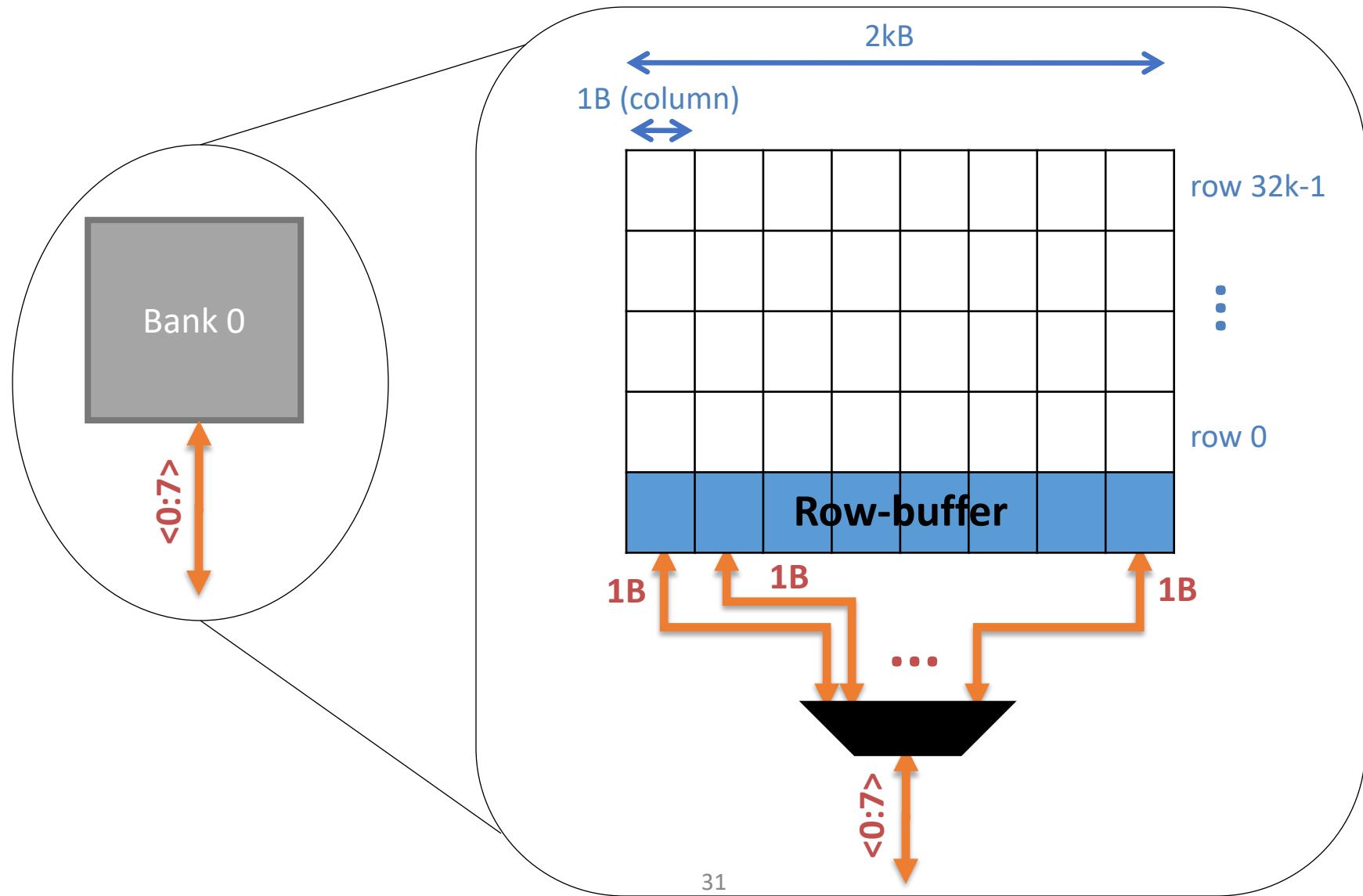
# Breaking down a Rank



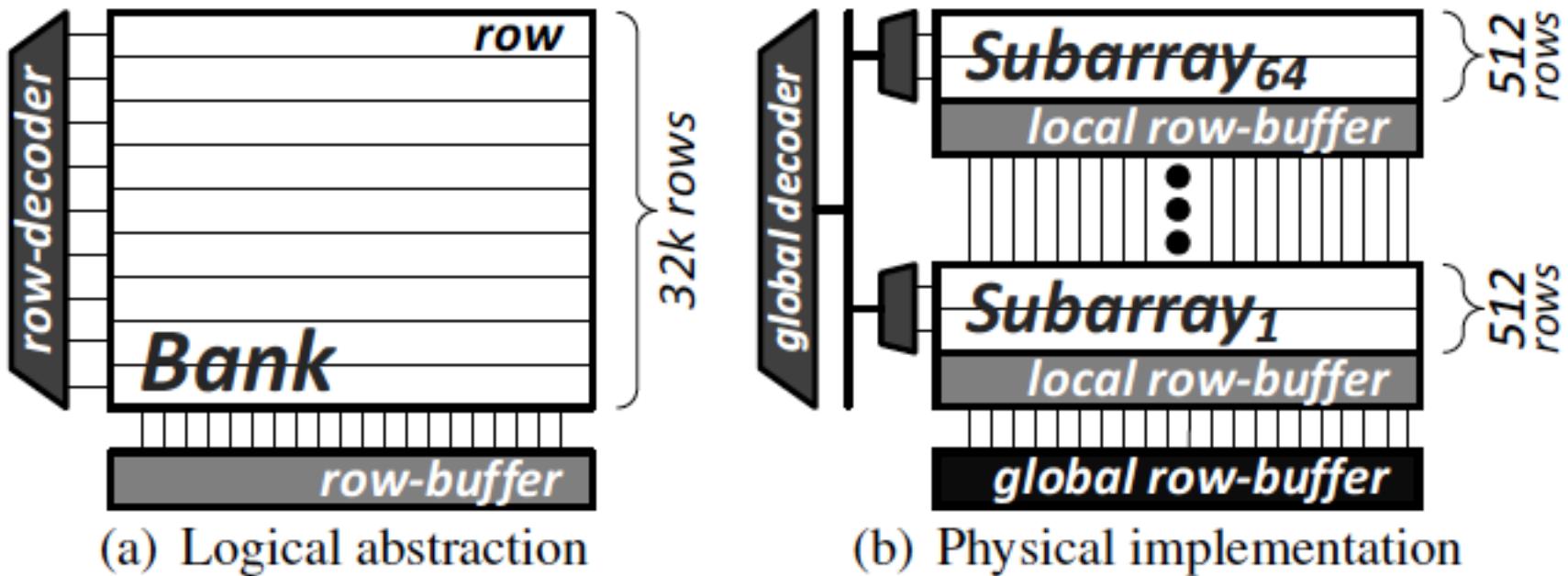
# Breaking down a Chip



# Breaking down a Bank

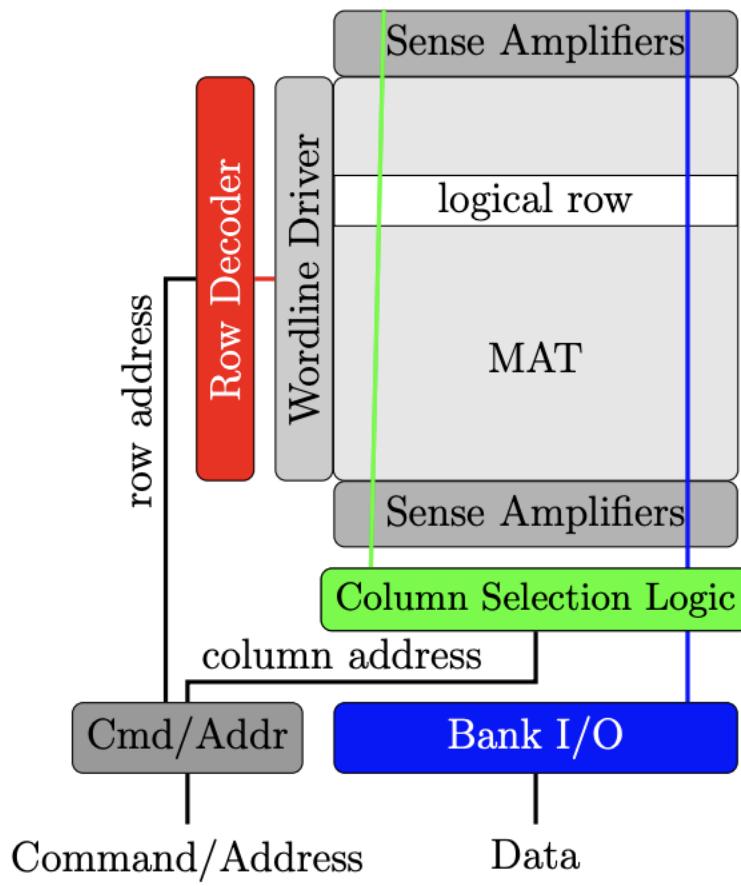


# A DRAM Bank Internally Has Sub-Banks

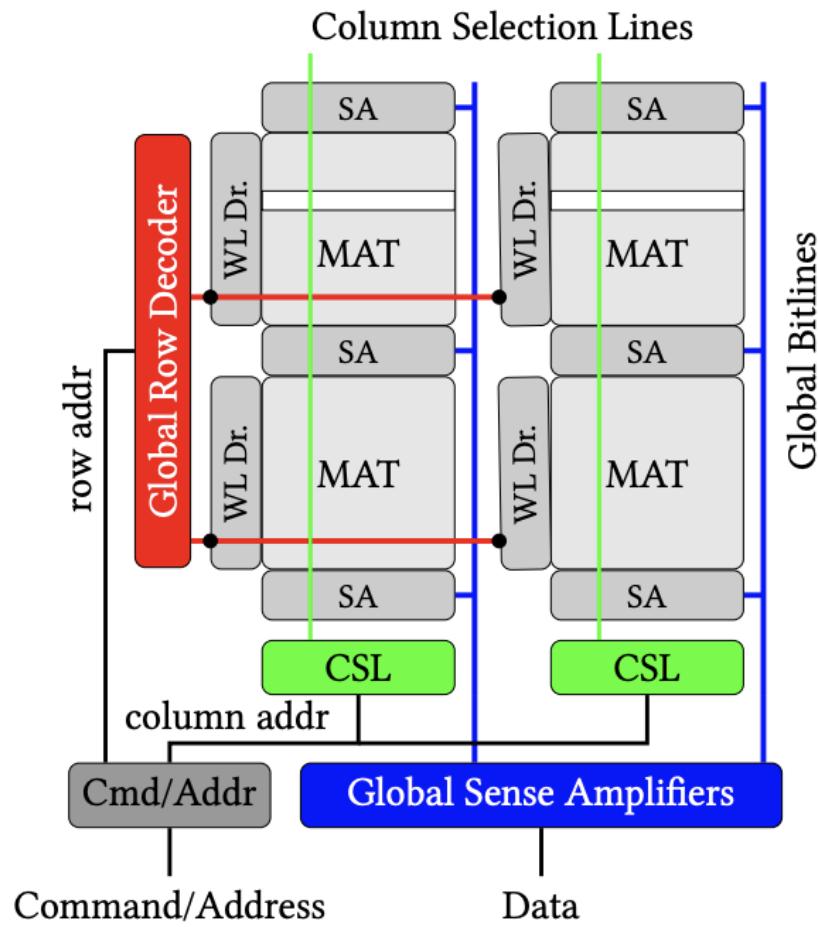


**Figure 1.** DRAM bank organization

# Another View of a DRAM Bank



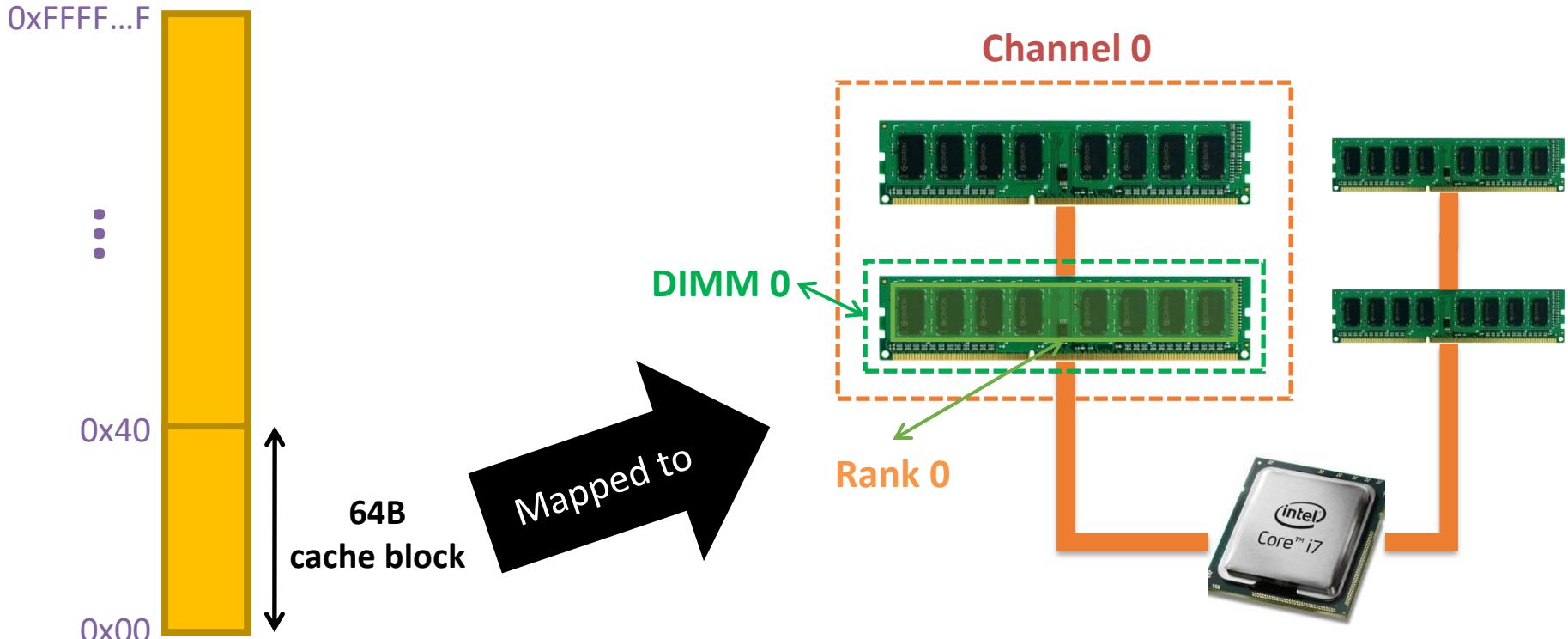
**Logical Abstraction**



**Physical View**

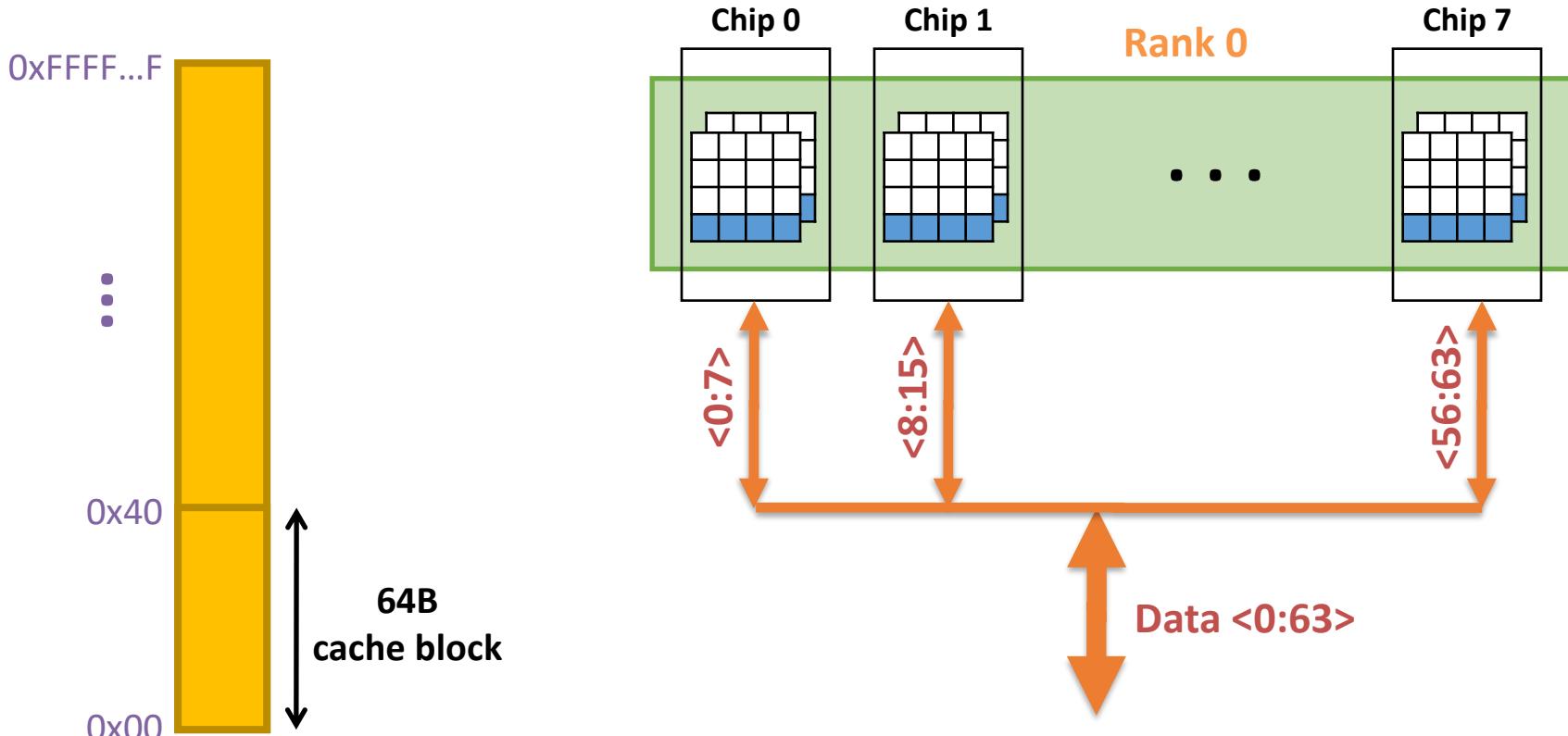
# Example: Transferring a cache block

Physical memory space



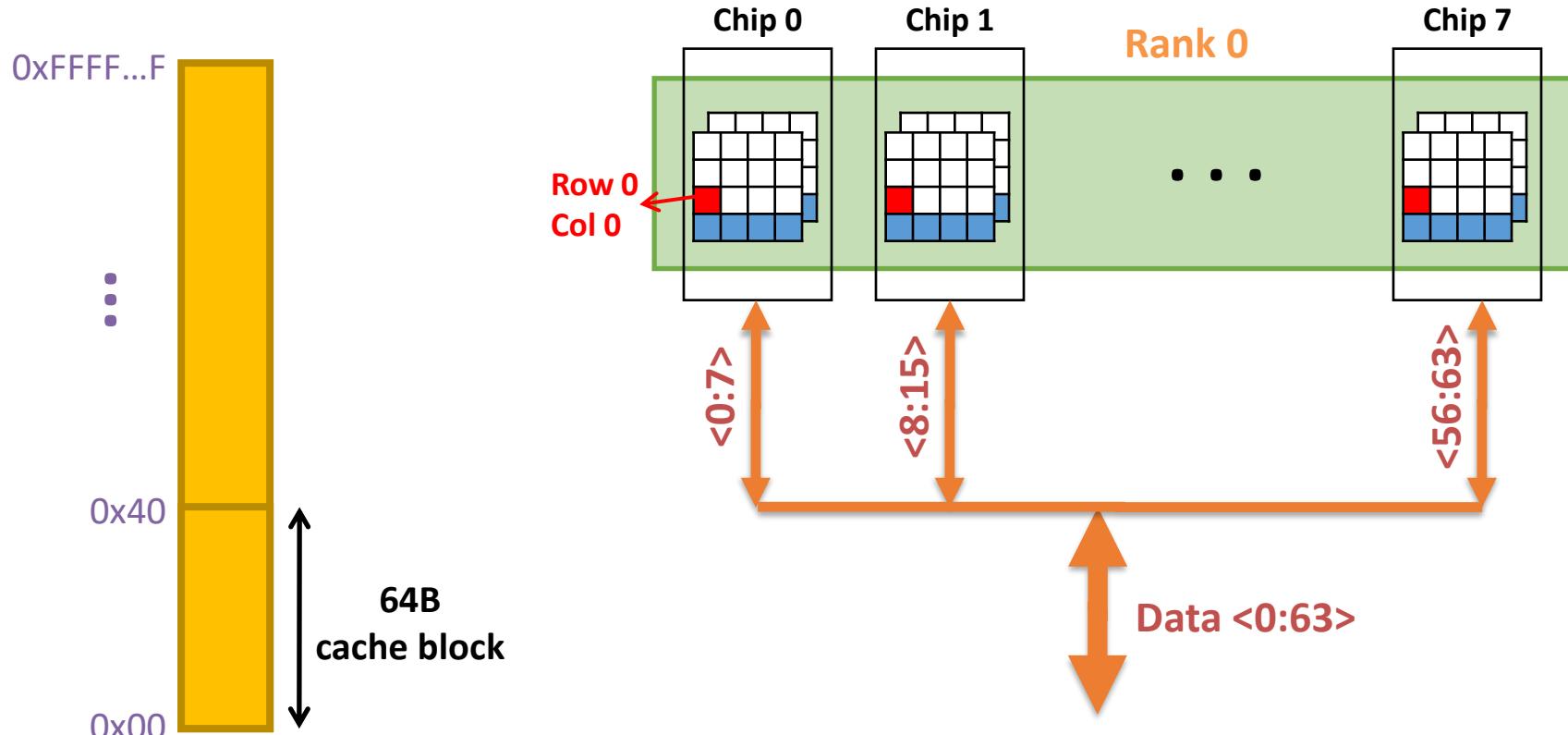
# Example: Transferring a cache block

Physical memory space



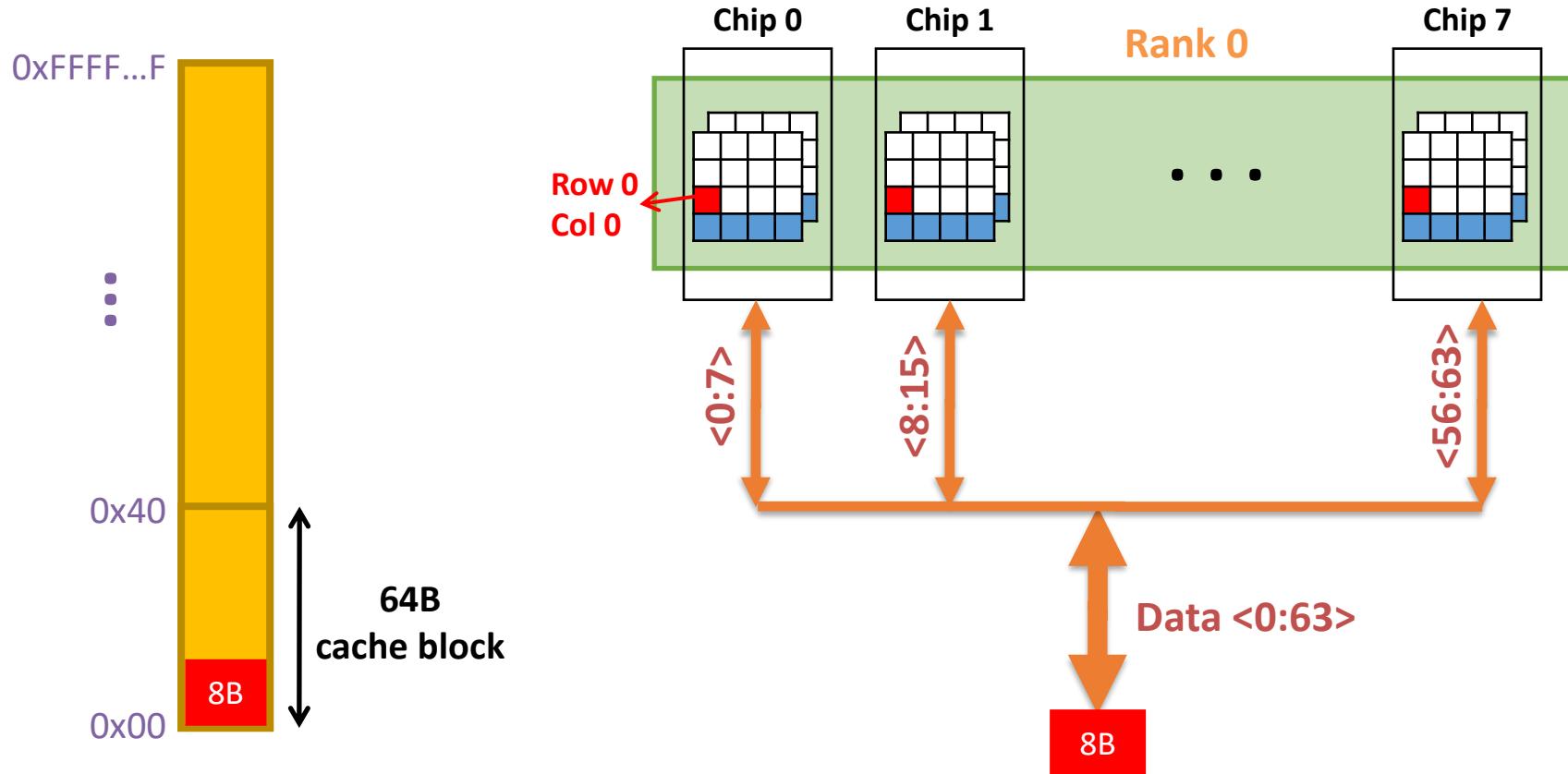
# Example: Transferring a cache block

Physical memory space



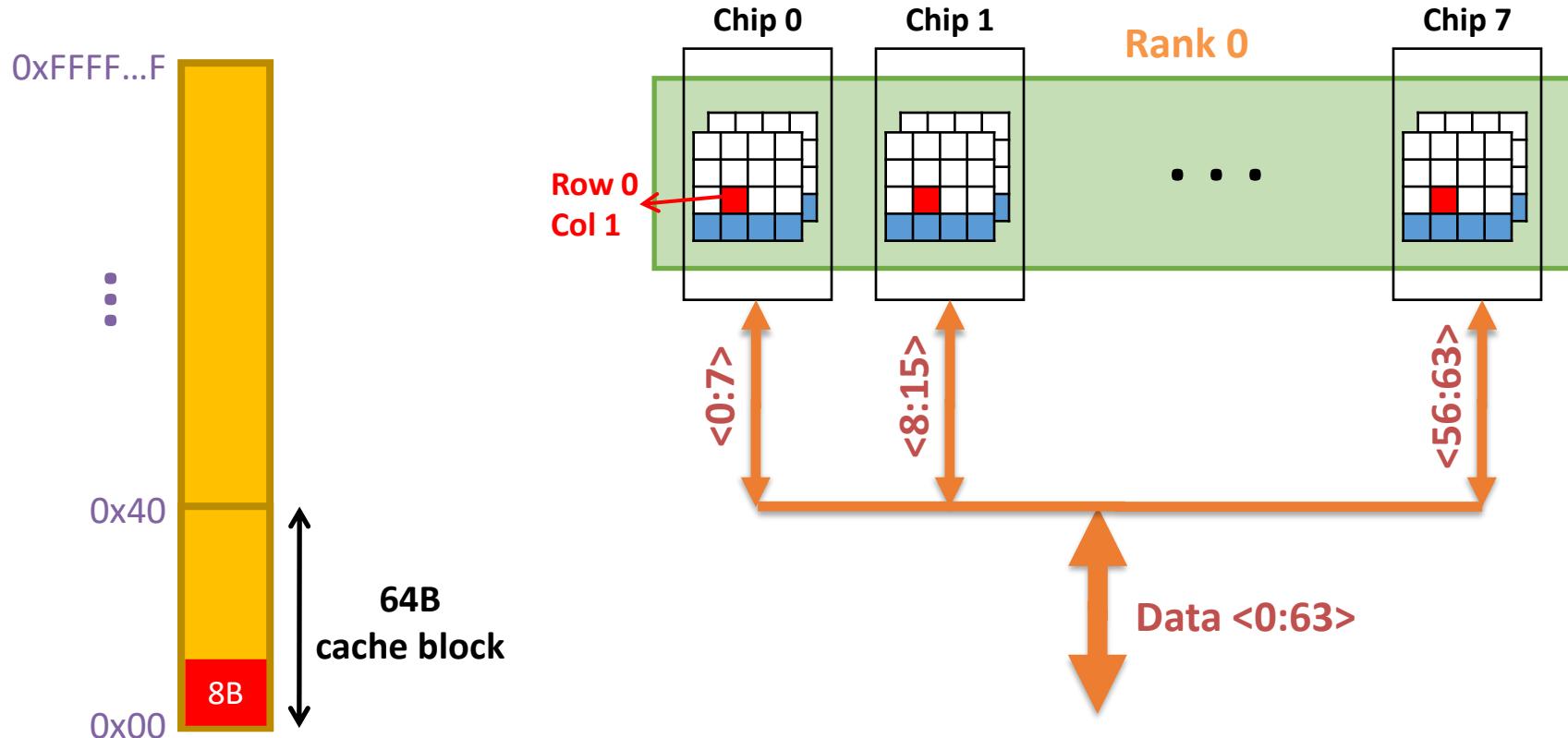
# Example: Transferring a cache block

Physical memory space



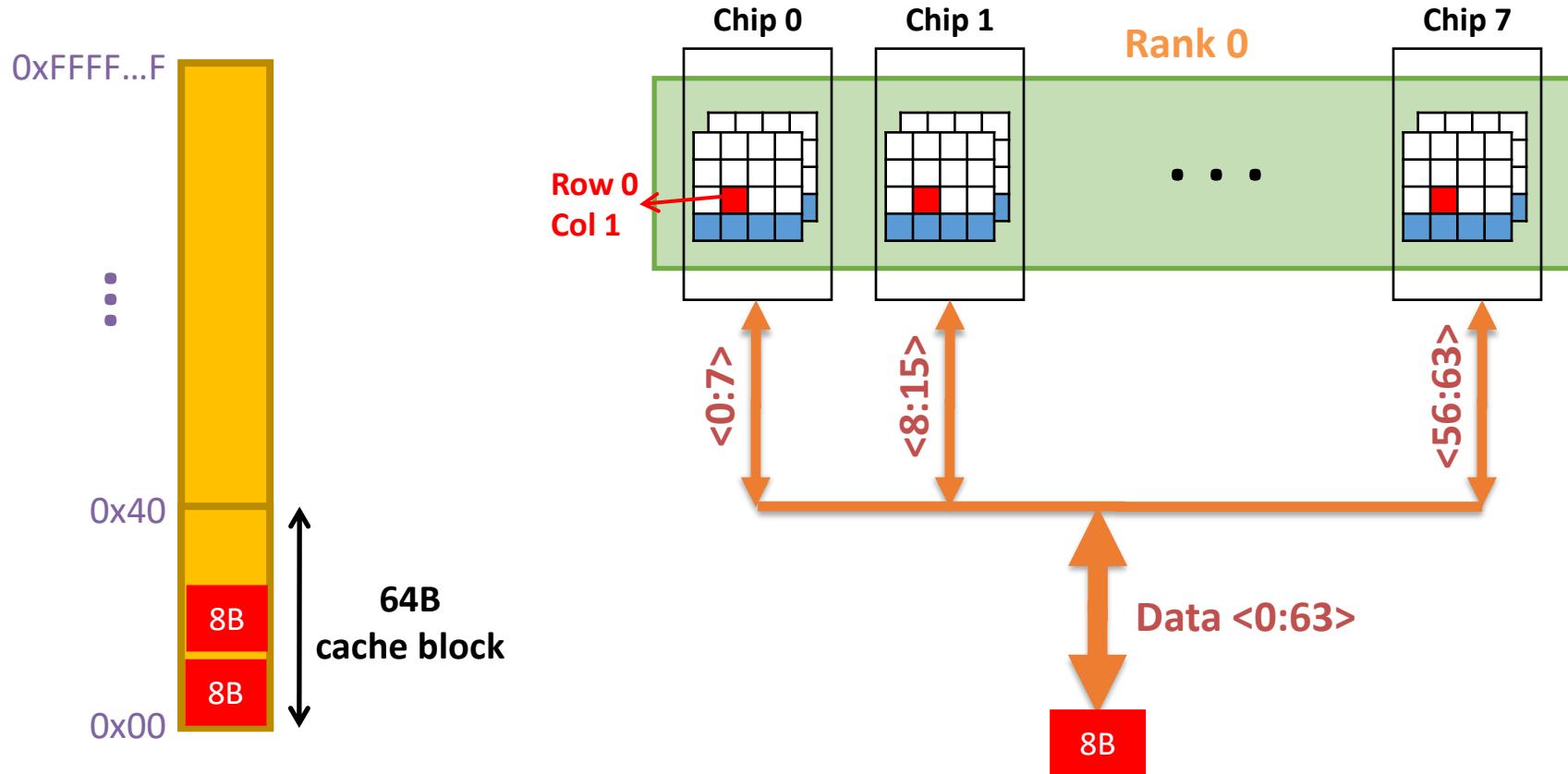
# Example: Transferring a cache block

Physical memory space



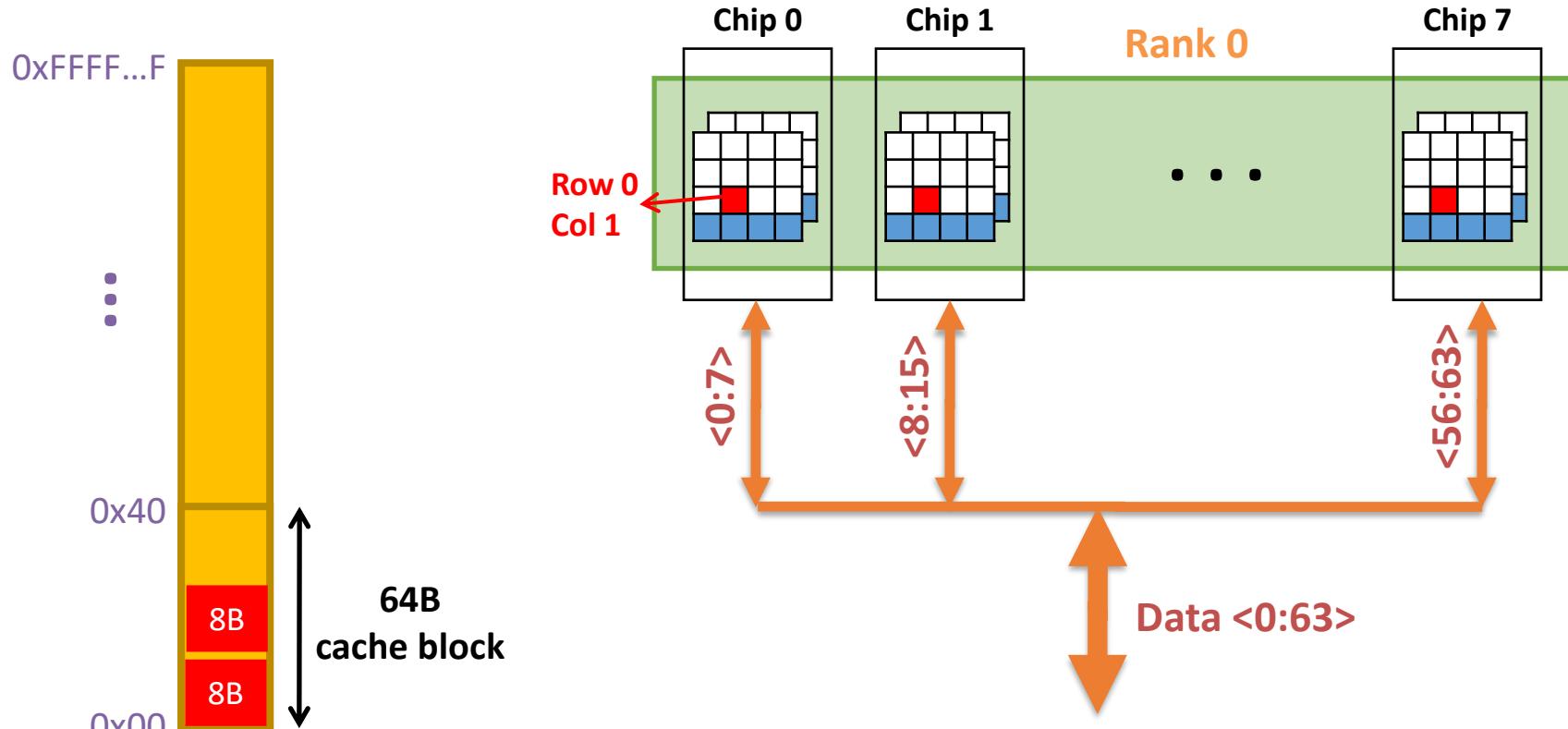
# Example: Transferring a cache block

Physical memory space



# Example: Transferring a cache block

Physical memory space



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

During the process, 8 columns are read sequentially.

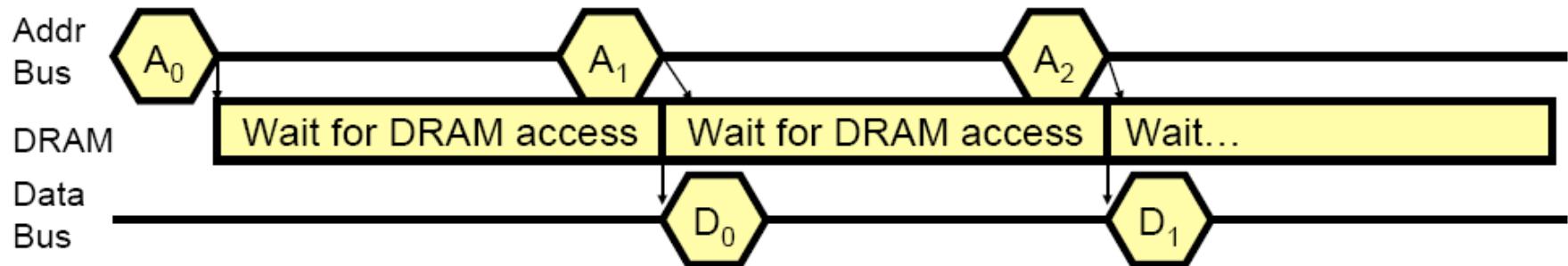
# Latency Components: Basic DRAM Operation

- 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

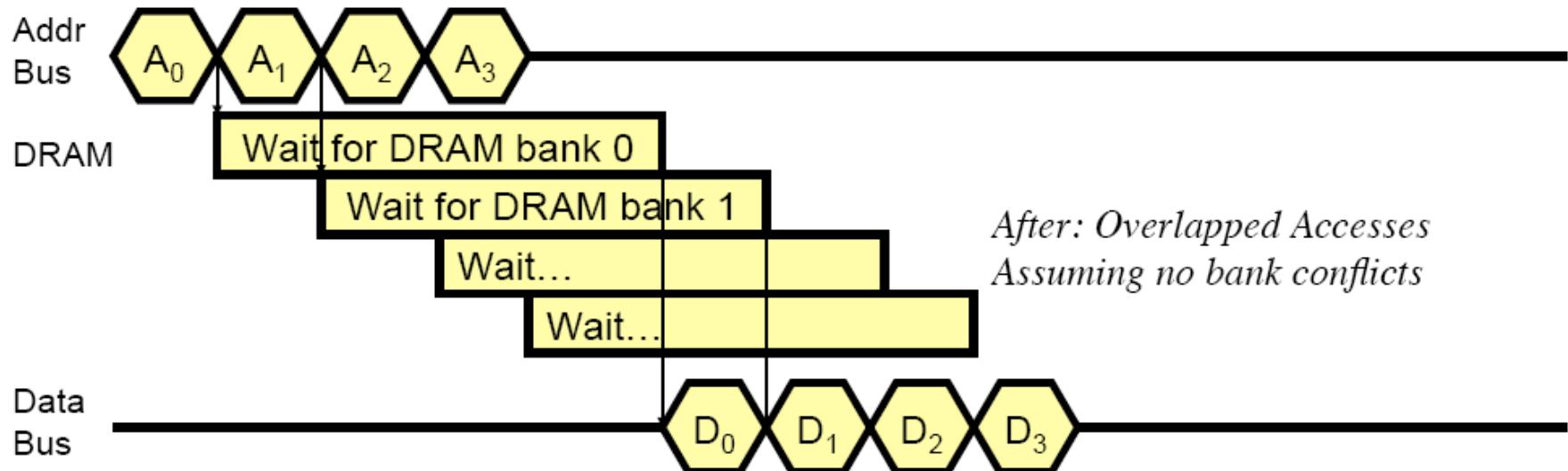
# Multiple Banks (Interleaving) and Channels

- 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/Channels Help



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



*After: Overlapped Accesses  
Assuming no bank conflicts*

# Multiple Channels

- Advantages
  - Increased bandwidth
  - Multiple concurrent accesses (if independent channels)
- Disadvantages
  - Higher cost than a single channel
    - More board wires
    - More pins (if on-chip memory controller)

# 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)
---------------	---------------	------------------	----------------------

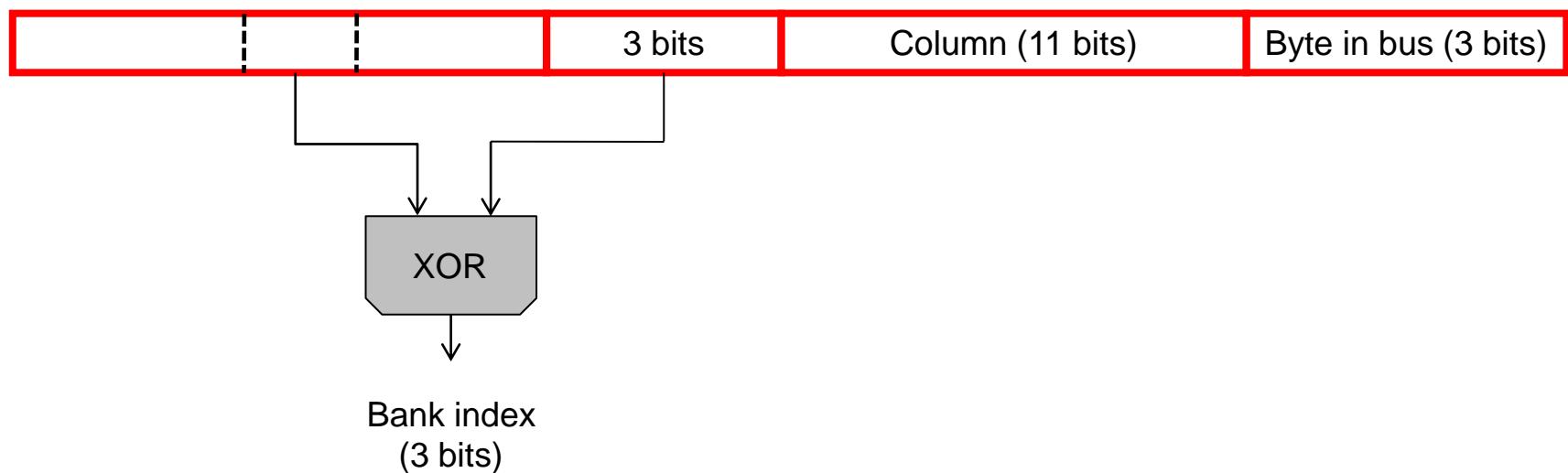
- Accesses to consecutive cache blocks serviced in a pipelined manner
- Cache block interleaving
  - Consecutive cache block addresses in consecutive banks
  - 64 byte cache blocks

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



# 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)

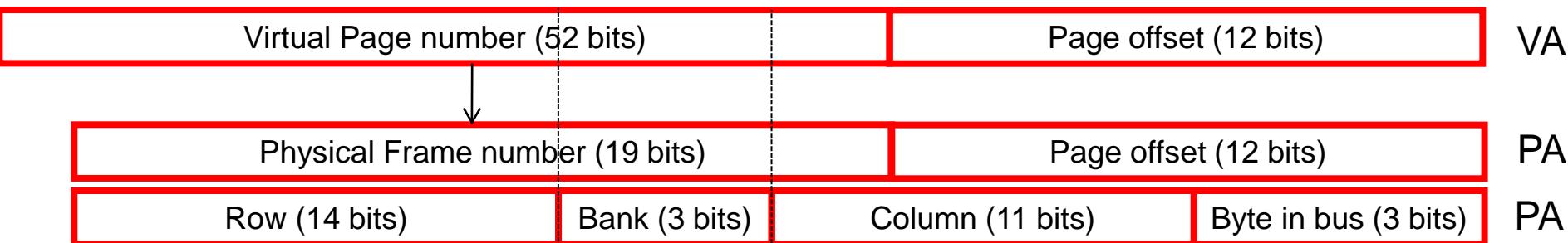
- Where are consecutive cache blocks?

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)	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		

47

# Interaction with Virtual → Physical Mapping

- Operating System influences where an address maps to in DRAM

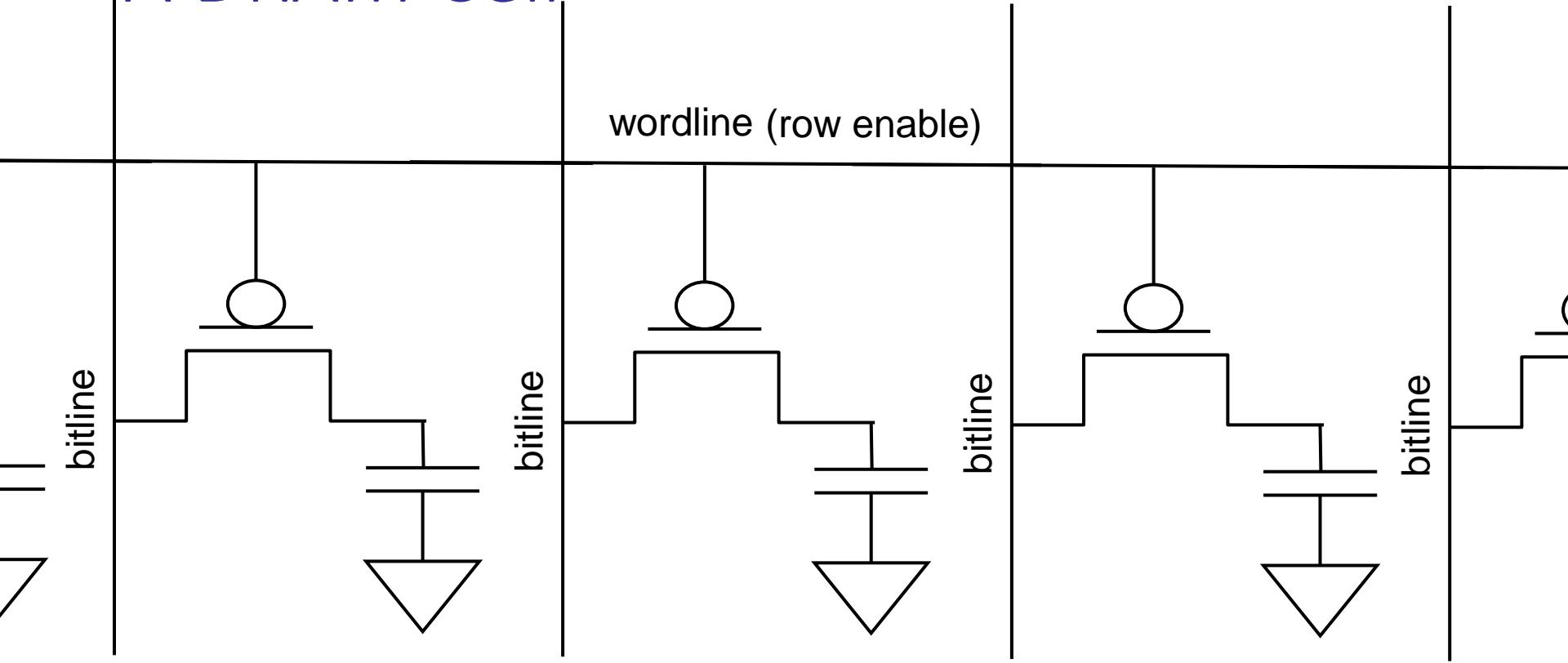


- Operating system can influence which bank/channel/rank a virtual page is mapped to.
- It can perform page coloring to
  - Minimize bank conflicts
  - Minimize inter-application interference **[Muralidhara + MICRO'11]**

# DRAM Refresh

(本節內容改自 Prof. Onur Mutlu, “Computer Architecture,” 8<sup>th</sup> Lecture, Fall 2023 課程講義)

# A DRAM Cell

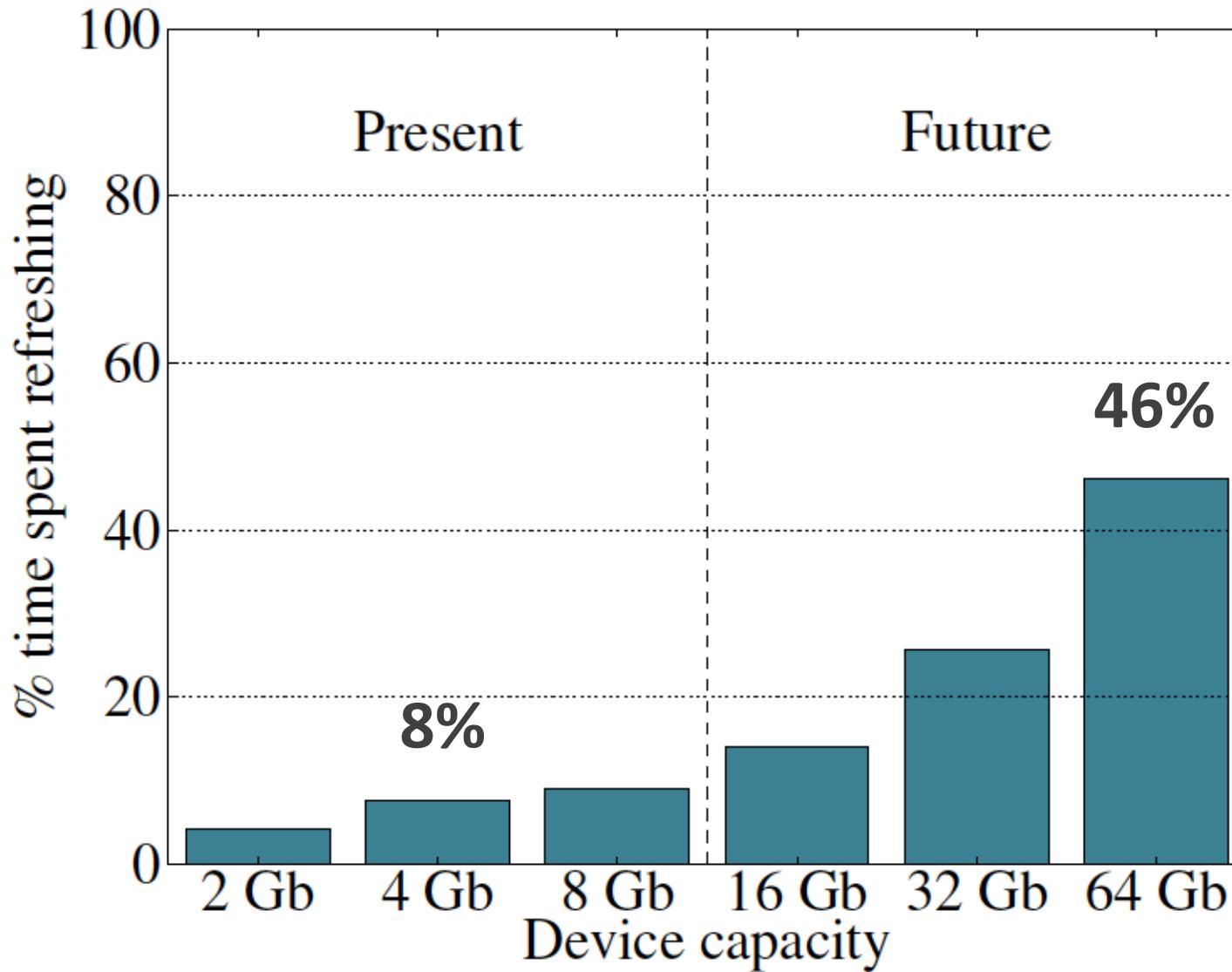


- A DRAM cell consists of a capacitor and an access transistor
- It stores data in terms of charge in the capacitor
- A DRAM chip consists of (10s of 1000s of) rows of such cells

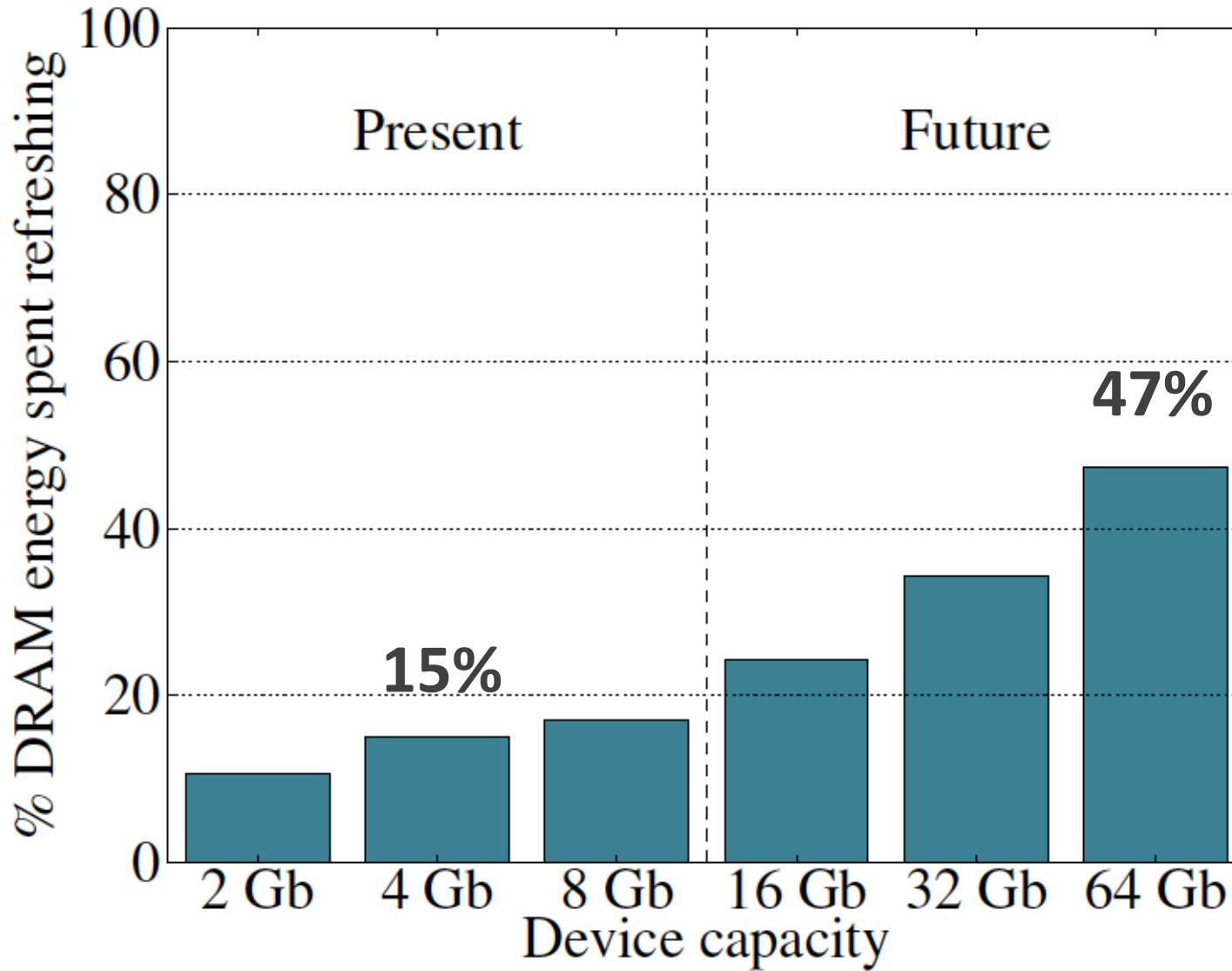
# DRAM Refresh

- DRAM capacitor charge leaks over time
- The memory controller needs to refresh each row periodically to restore charge
  - Activate 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 capacity scaling

# Refresh Overhead: Performance



# Refresh Overhead: Energy

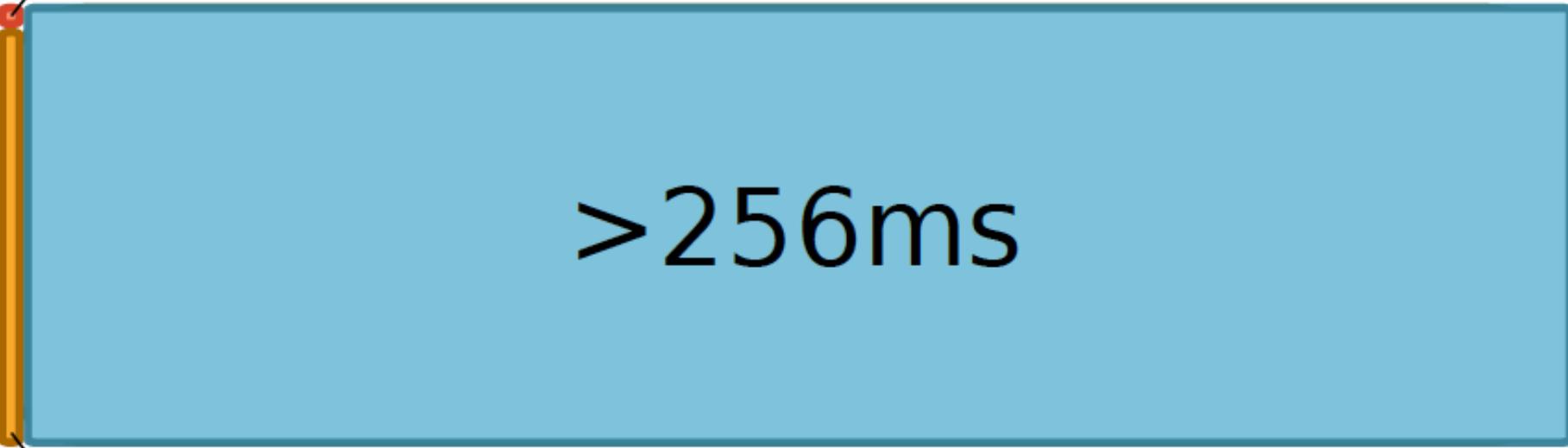


# How Do We Solve the Problem?

- Observation: All DRAM rows are refreshed every 64ms.
- Critical thinking: Do we need to refresh all rows every 64ms?
- What if we knew what happened underneath (in DRAM cells) and exposed that information to upper layers?

# Underneath: Retention Time Profile of DRAM

64-128ms



>256ms

128-256ms

# Aside: Why Do We Have Such a Profile?

- Answer: Manufacturing is not perfect
- Not all DRAM cells are exactly the same
- Some cells are more leaky than others
- This is called **Manufacturing Process Variation**

# Opportunity: Taking Advantage of This Profile

- Assume we know the retention time of each row exactly
- What can we do with this information?
- Who do we expose this information to?
- How much information do we expose?
  - Affects hardware/software overhead, power verification complexity, cost
- How do we determine this profile information?
  - Also, who determines it?

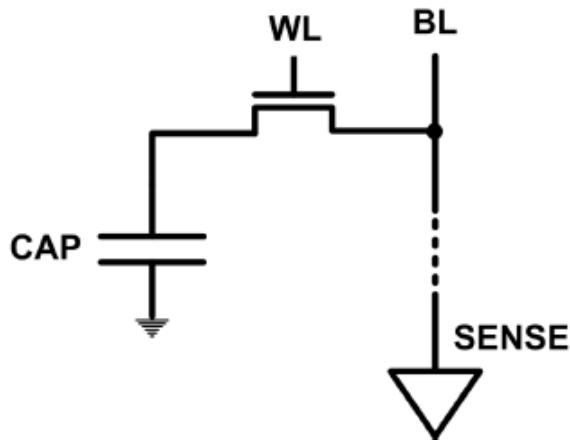
Problem
Algorithm
Program/Language
Runtime System (VM, OS, MM)
ISA (Architecture)
Microarchitecture
Logic
Devices
Electrons

# DRAM Reliability Issues

(本節內容改自 Prof. Onur Mutlu, “Computer Architecture,” 6<sup>th</sup> Lecture, Fall 2023 課程講義)

# The DRAM Scaling Problem

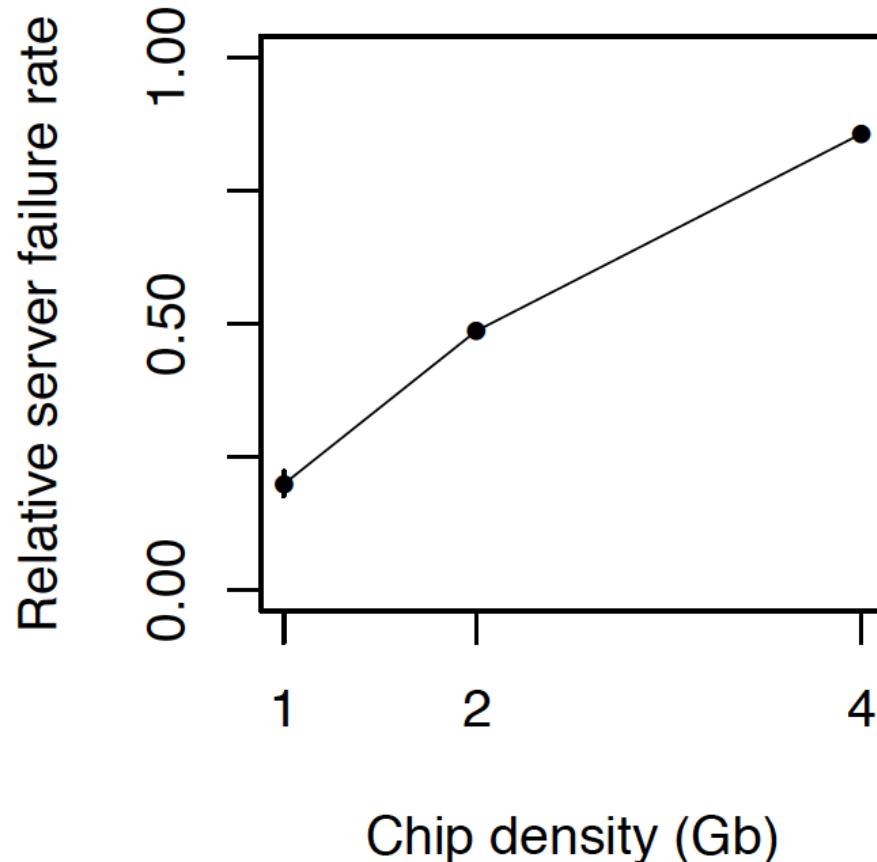
- DRAM stores charge in a capacitor (charge-based memory)
  - Capacitor must be large enough for reliable sensing
  - Access transistor should be large enough for low leakage and high retention time
  - Scaling beyond 40-35nm (2013) is challenging [ITRS, 2009]



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

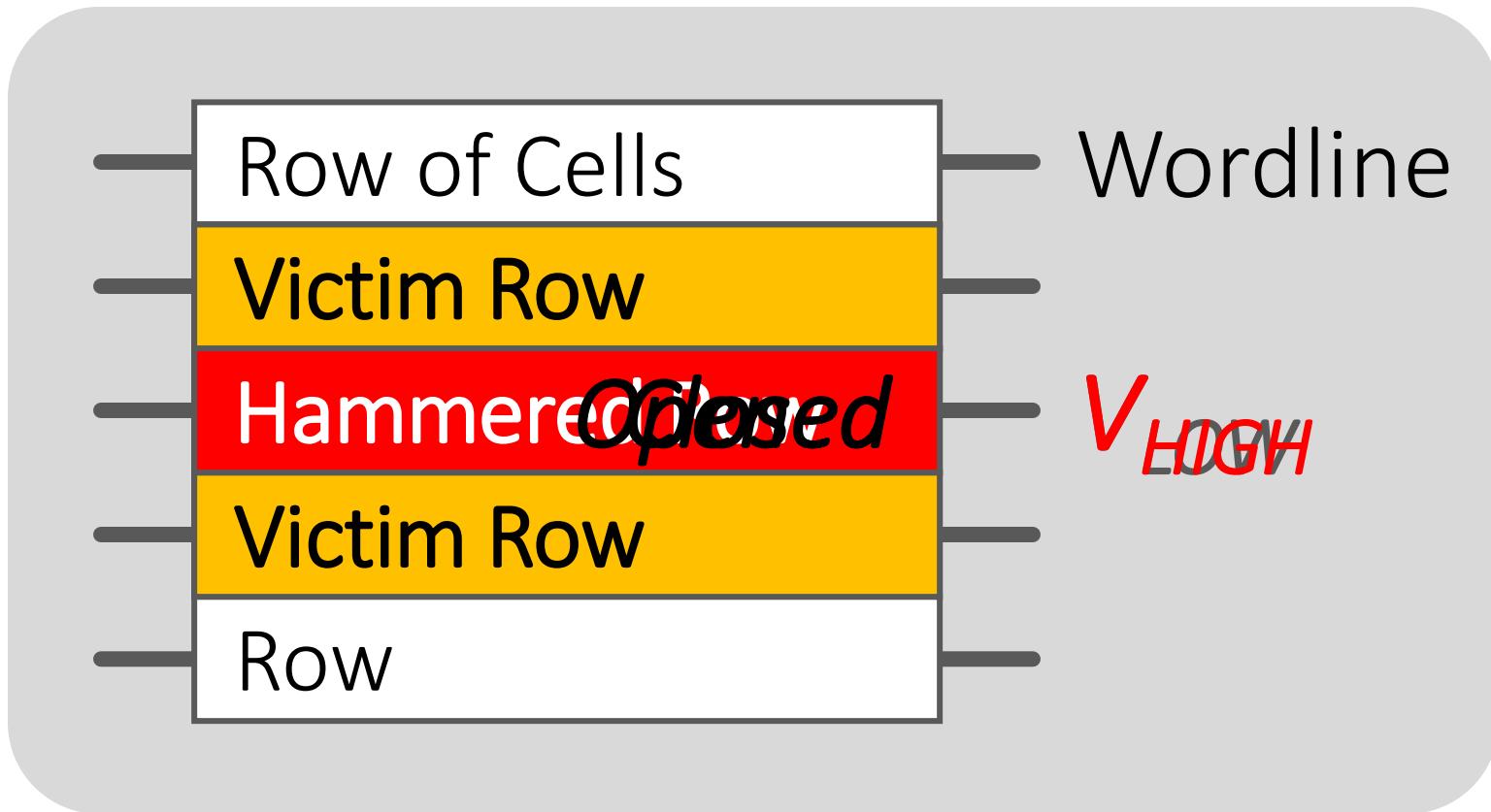
# As Memory Scales, It Becomes Unreliable

- Data from all of Facebook's servers worldwide
- Meza+, “[Revisiting Memory Errors in Large-Scale Production Data Centers](#),” DSN’15.



*Intuition:  
quadratic  
increase  
in  
capacity*

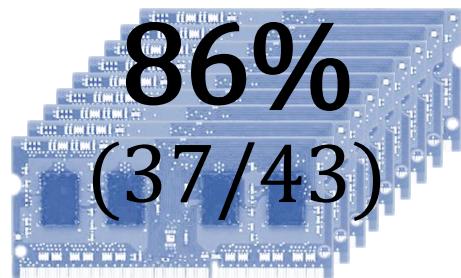
# Modern DRAM is Prone to Disturbance Errors



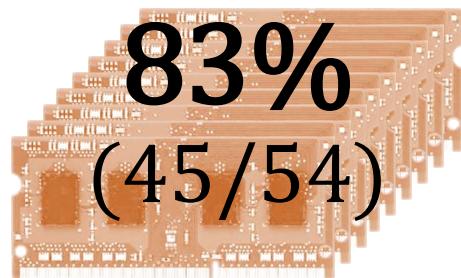
Repeatedly reading a row enough times (before memory gets refreshed) induces **disturbance errors** in adjacent rows in **most real DRAM chips you can buy today**

# Most DRAM Modules Are Vulnerable

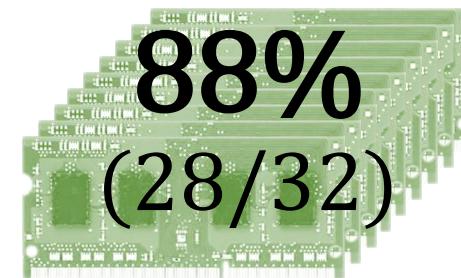
A company



B company



C company

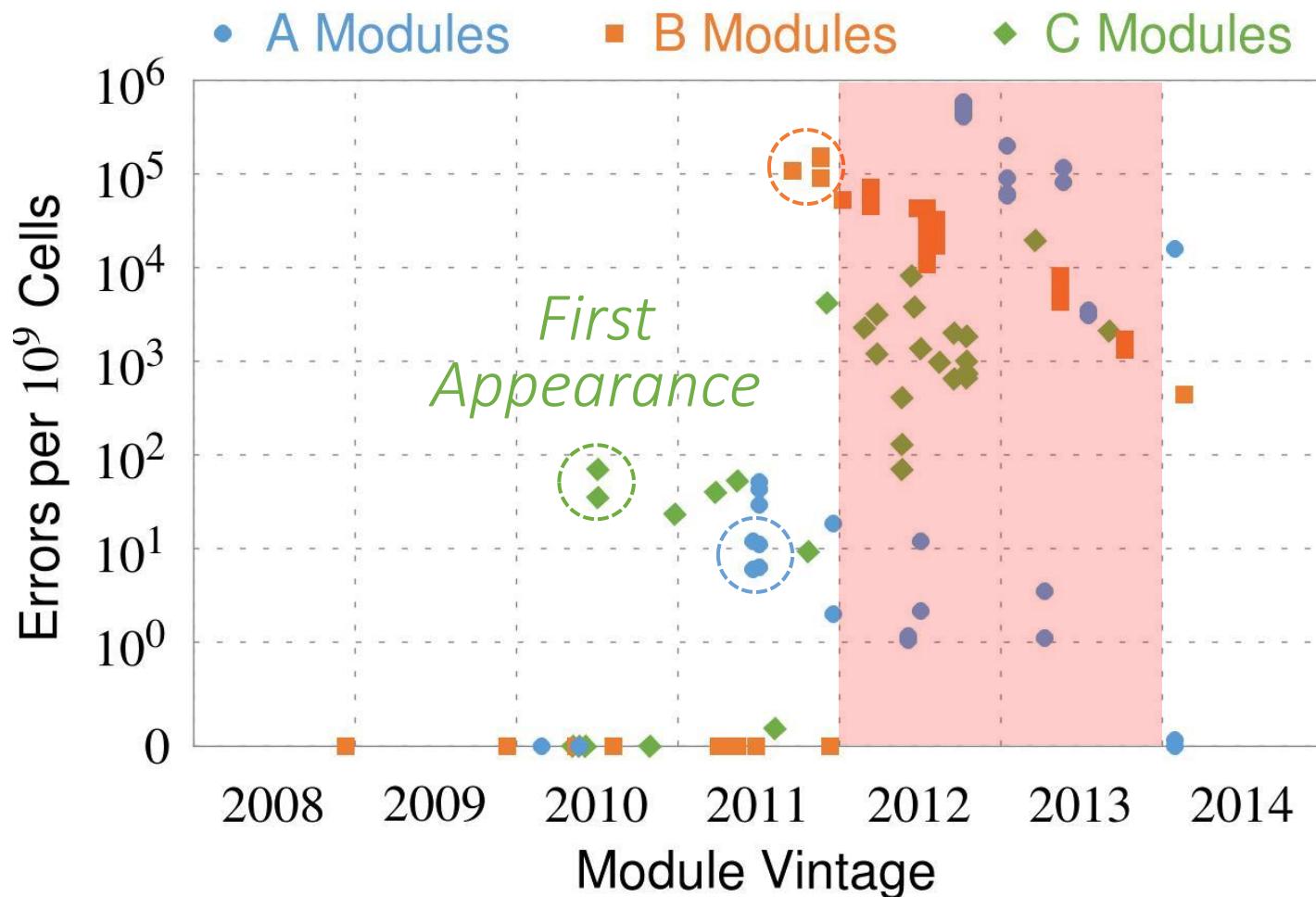


Up to  
 $1.0 \times 10^7$   
errors

Up to  
 $2.7 \times 10^6$   
errors

Up to  
 $3.3 \times 10^5$   
errors

# Recent DRAM Is More Vulnerable



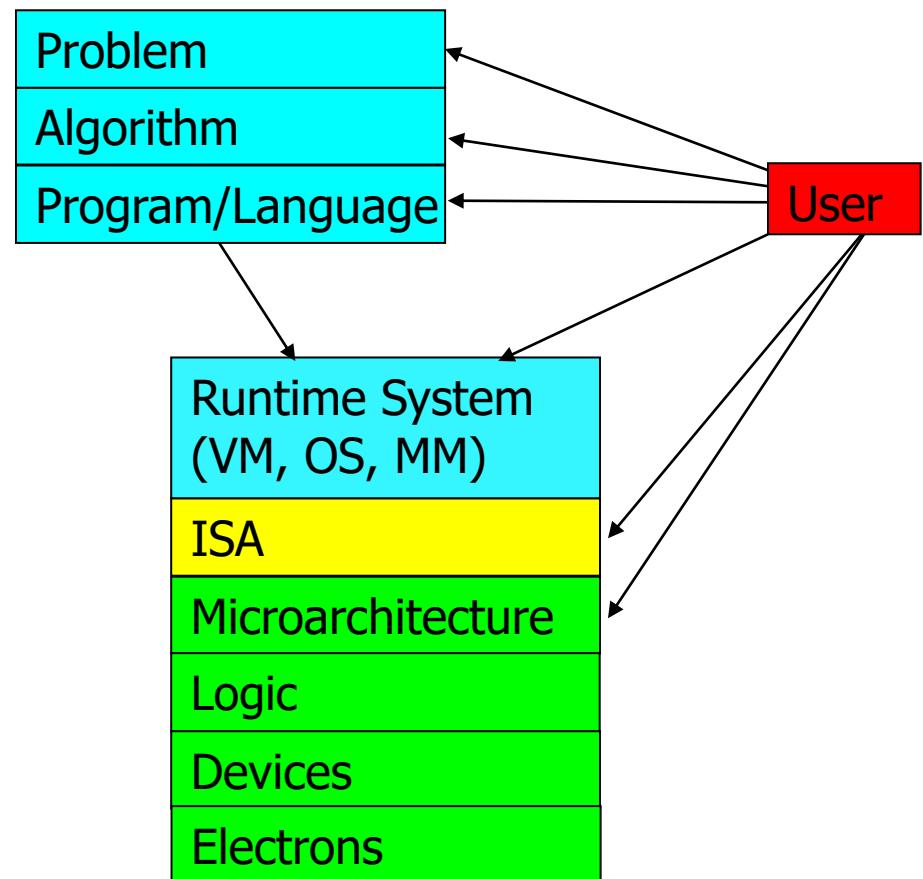
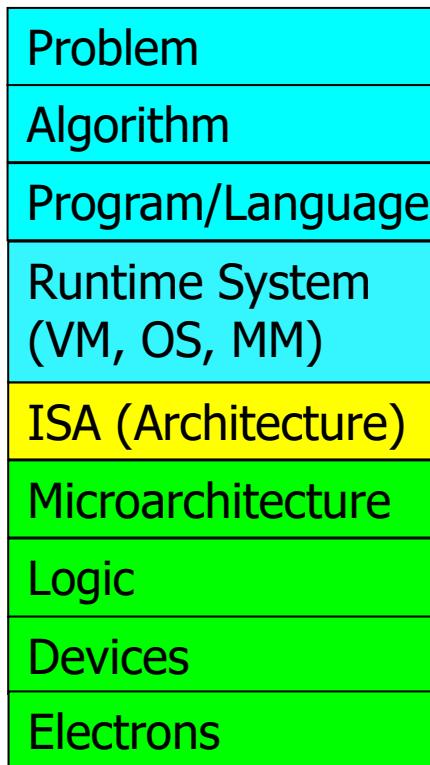
*All modules from 2012–2013 are vulnerable*

# Why Is This Happening?

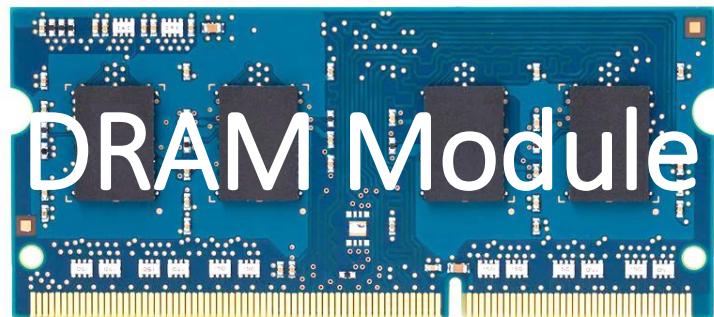
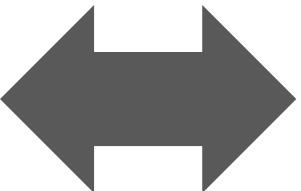
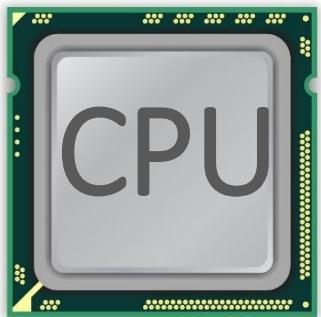
- DRAM cells are too close to each other!
  - They are not electrically isolated from each other
- Access to one cell affects the value in nearby cells
  - due to **electrical interference** between
    - the cells
    - wires used for accessing the cells
  - Also called cell-to-cell coupling/interference
- Example: When we activate (apply high voltage) to a row, an adjacent row gets slightly activated as well
  - Vulnerable cells in that slightly-activated row lose a little bit of charge
  - If RowHammer happens enough times, charge in such cells gets drained

# Higher-Level Implications

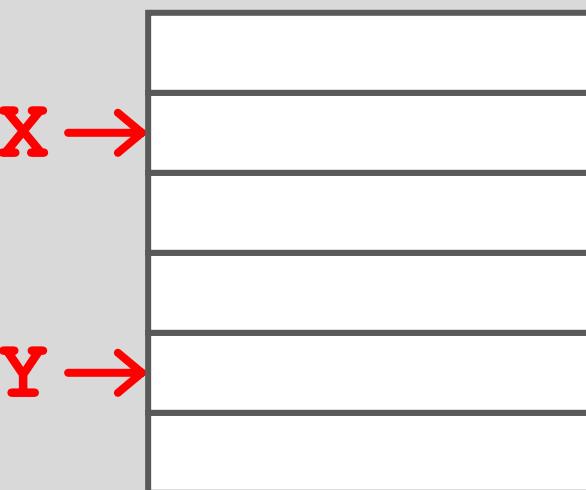
- This simple circuit level failure mechanism has enormous implications on upper layers of the transformation hierarchy



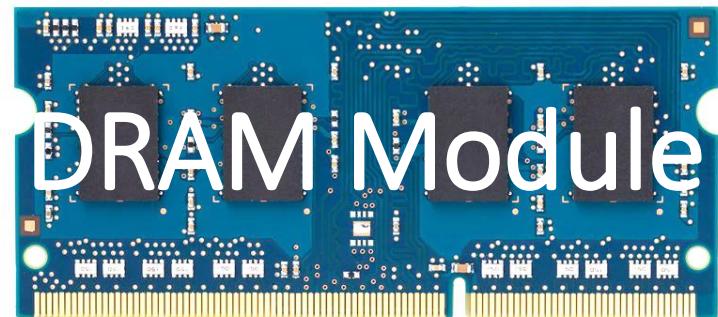
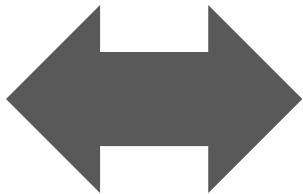
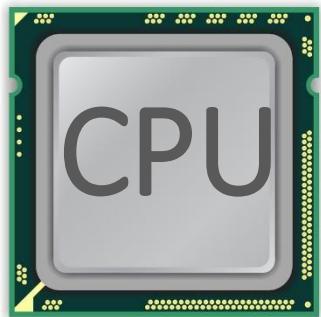
# A Simple Program Can Induce Many Errors



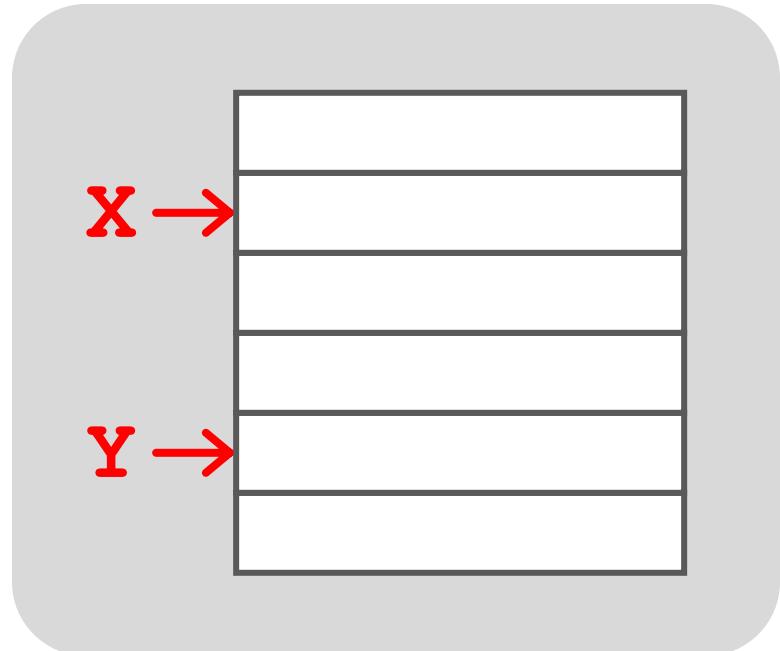
```
loop:  
    mov (%X), %eax  
    mov (%Y), %ebx  
    clflush (%X)  
    clflush (%Y)  
    mfence  
    jmp loop
```



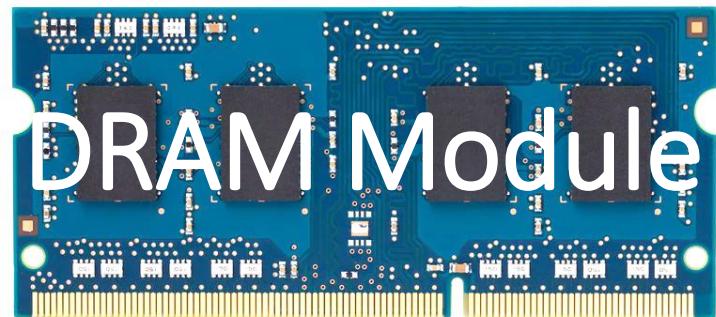
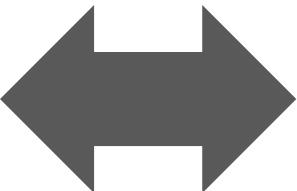
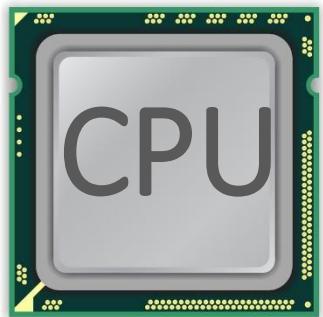
# A Simple Program Can Induce Many Errors



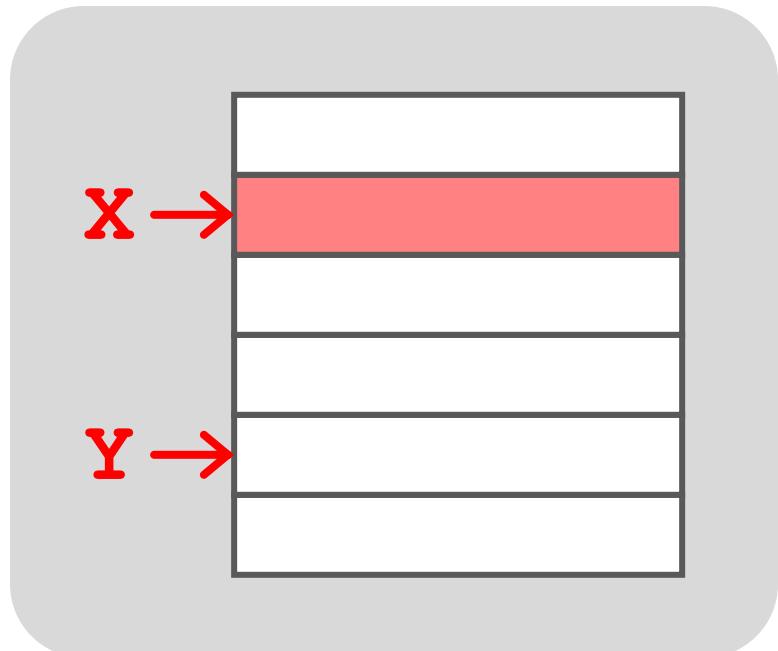
1. Avoid *cache hits*
  - Flush **X** from cache
2. Avoid *row hits* to **X**
  - Read **Y** in another row



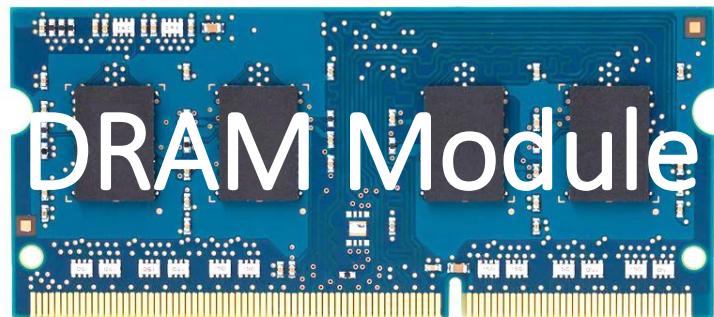
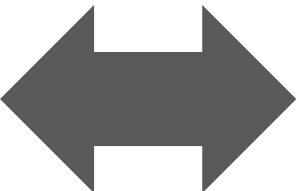
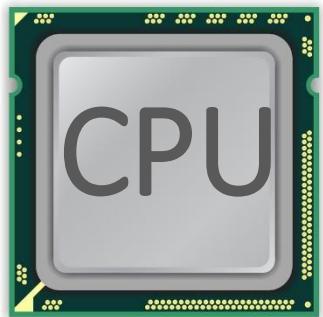
# A Simple Program Can Induce Many Errors



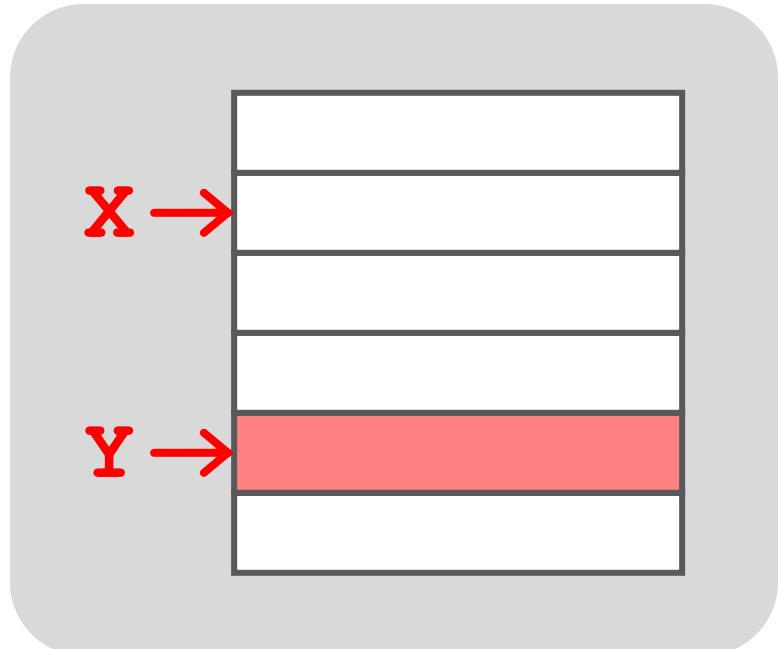
```
loop:  
    mov (%X), %eax  
    mov (%Y), %ebx  
    clflush (%X)  
    clflush (%Y)  
    mfence  
    jmp loop
```



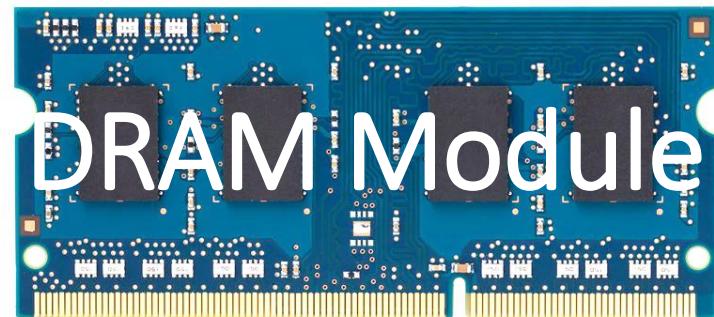
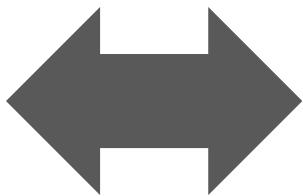
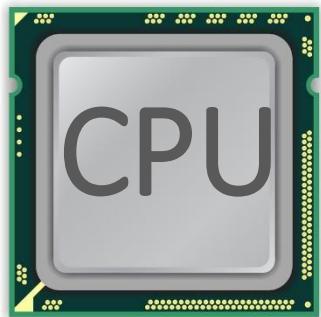
# A Simple Program Can Induce Many Errors



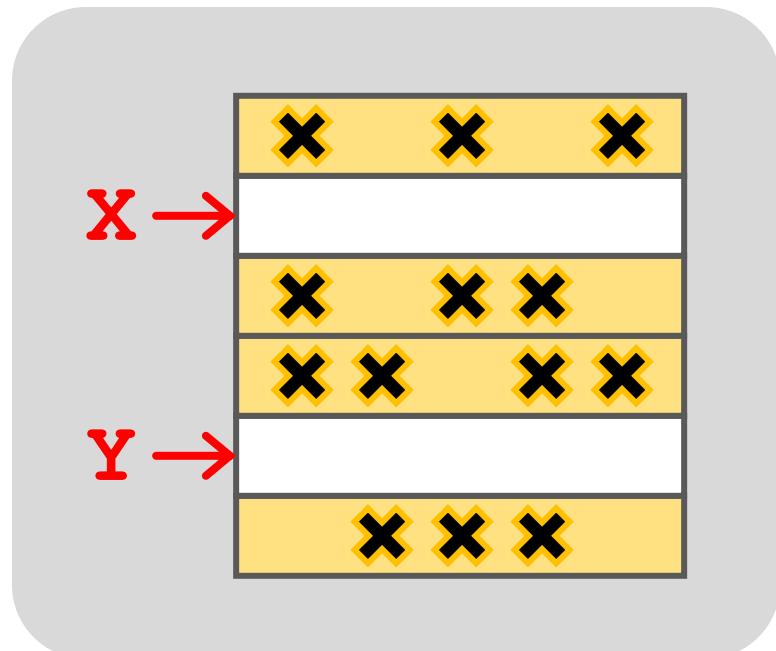
```
loop:  
    mov (%X), %eax  
    mov (%Y), %ebx  
    clflush (%X)  
    clflush (%Y)  
    mfence  
    jmp loop
```



# A Simple Program Can Induce Many Errors



```
loop:  
    mov (%X), %eax  
    mov (%Y), %ebx  
    clflush (%X)  
    clflush (%Y)  
    mfence  
    jmp loop
```



# Observed Errors in Real Systems

CPU Architecture	Errors	Access-Rate
Intel Haswell (2013)	22.9K	12.3M/sec
Intel Ivy Bridge (2012)	20.7K	11.7M/sec
Intel Sandy Bridge (2011)	16.1K	11.6M/sec
AMD Piledriver (2012)	59	6.1M/sec

**A real reliability & security issue**

# Security Implications



# Security Implications



It's like breaking into an apartment by repeatedly slamming a neighbor's door until the vibrations open the door you were after

# More Security Implications (I)

**"We can gain unrestricted access to systems of website visitors."**

www.iaik.tugraz.at ■

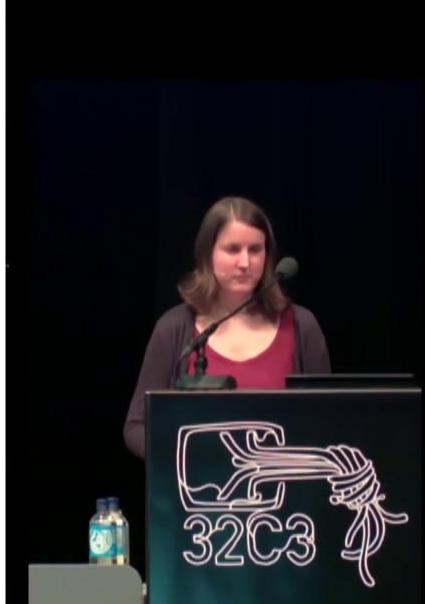
Not there yet, but ...



ROOT privileges for web apps!

29

Daniel Gruss (@lavados), Clémentine Maurice (@BloodyTangerine),  
December 28, 2015 — 32c3, Hamburg, Germany



GATED  
COMMUNITIES

Rowhammer.js: A Remote Software-Induced Fault Attack in JavaScript (DIMVA'16)

74

Source: <https://lab.dsst.io/32c3-slides/7197.html>

# More Security Implications (II)

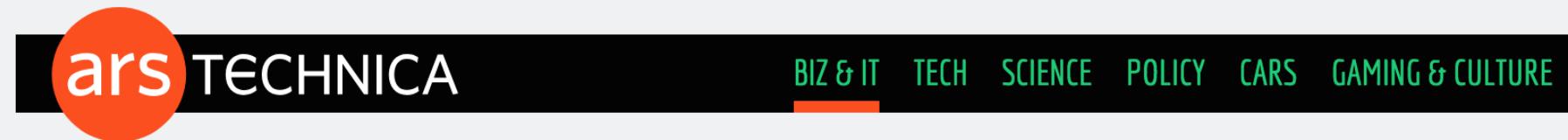
**"Can gain control of a smart phone deterministically"**



Drammer: Deterministic Rowhammer  
Attacks on Mobile Platforms, CCS'16

# More Security Implications (III)

- Using an integrated GPU in a mobile system to remotely escalate privilege via the WebGL interface. [IEEE S&P](#)



"GRAND PWNING UNIT" —

## Drive-by Rowhammer attack uses GPU to compromise an Android phone

JavaScript based GLitch pwns browsers by flipping bits inside memory chips.

## Grand Pwning Unit: Accelerating Microarchitectural Attacks with the GPU

Pietro Frigo  
Vrije Universiteit  
Amsterdam  
p.frigo@vu.nl

Cristiano Giuffrida  
Vrije Universiteit  
Amsterdam  
giuffrida@cs.vu.nl

Herbert Bos  
Vrije Universiteit  
Amsterdam  
herbertb@cs.vu.nl

Kaveh Razavi  
Vrije Universiteit  
Amsterdam  
kaveh@cs.vu.nl

# Two Types of RowHammer Solutions

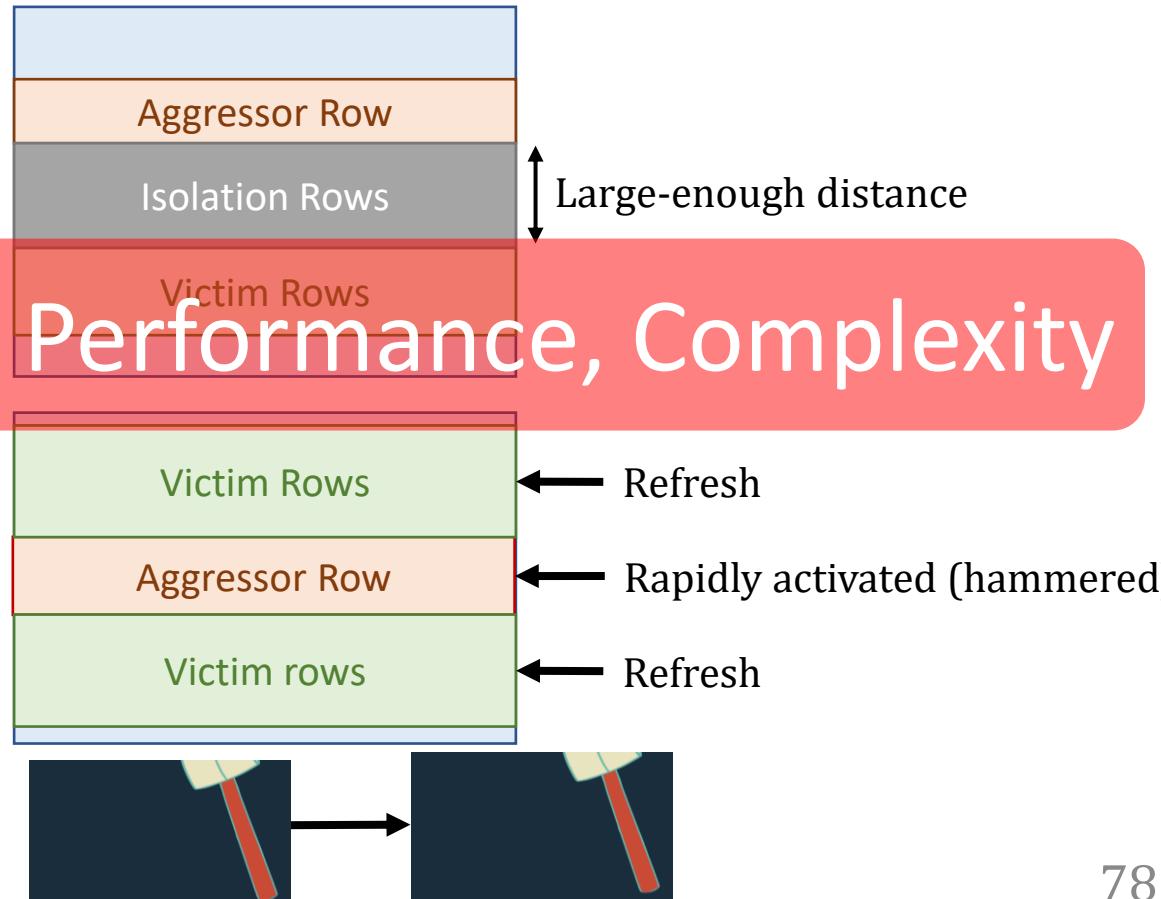
- Immediate
  - To protect the vulnerable DRAM chips in the field
  - Limited possibilities
- Longer-term
  - To protect future DRAM chips
  - Wider range of protection mechanisms

# RowHammer Solution Approaches

- More robust DRAM chips **and/or** error-correcting codes
- Increased refresh rate



- Physical isolation
- Cost, Power, Performance, Complexity  
Reactive refresh



- Proactive throttling

# Apple's Security Patch for RowHammer

- <https://support.apple.com/en-gb/HT204934>

Available for: OS X Mountain Lion v10.8.5, OS X Mavericks v10.9.5

Impact: A malicious application may induce memory corruption to escalate privileges

Description: A disturbance error, also known as Rowhammer, exists with some DDR3 RAM that could have led to memory corruption. This issue was mitigated by increasing memory refresh rates.

CVE-ID

CVE-2015-3693 : Mark Seaborn and Thomas Dullien of Google, working from original research by Yoongu Kim et al (2014)

HP, Lenovo, and many other vendors released similar patches

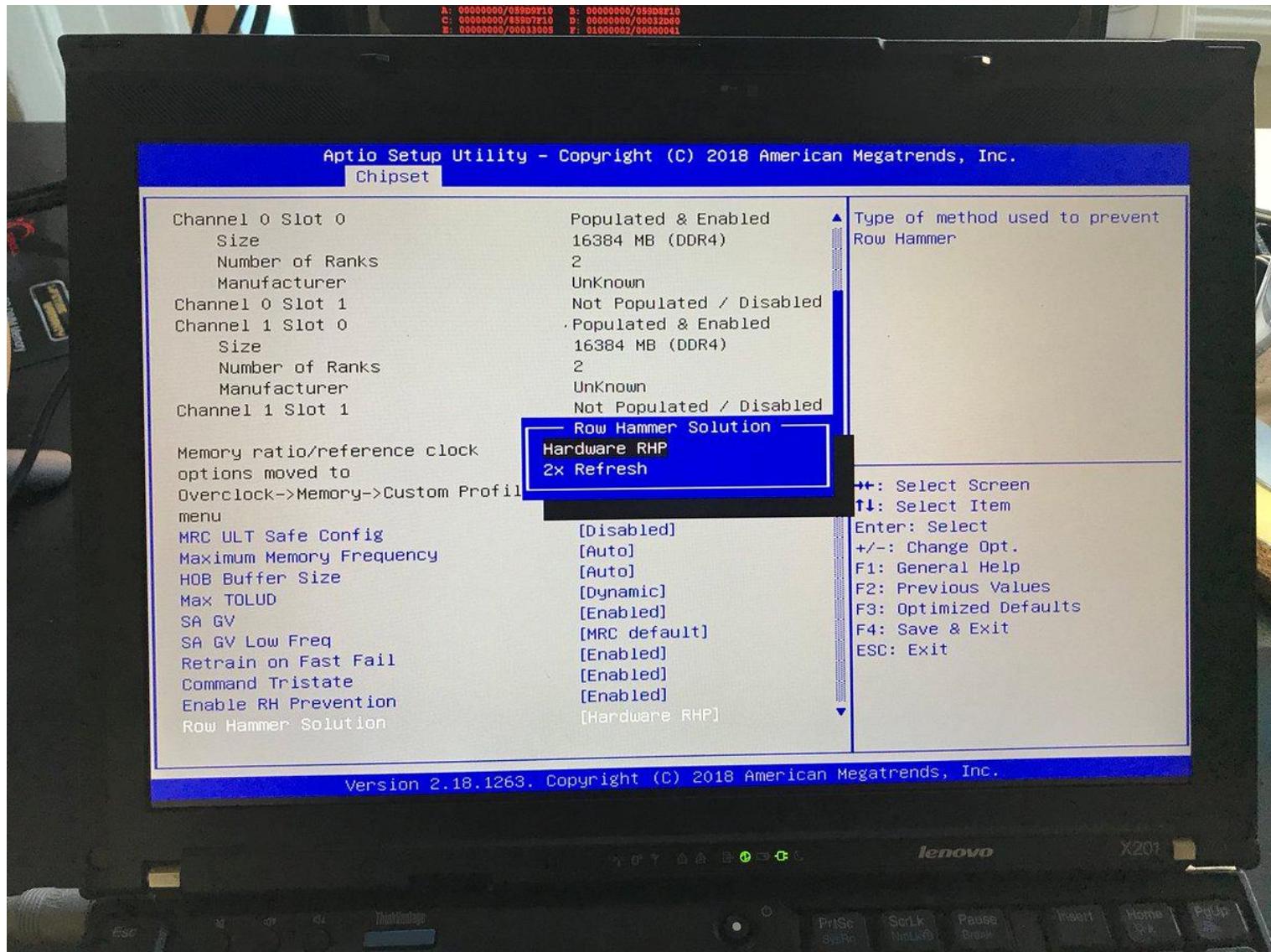
# PARA: Probabilistic Adjacent Row Activation

- Key Idea
  - After closing a row, we activate (i.e., refresh) one of its neighbors with a low probability:  $p = 0.005$
- Reliability Guarantee
  - When  $p=0.005$ , errors in one year:  $9.4 \times 10^{-14}$
  - By adjusting the value of  $p$ , we can vary the strength of protection against errors

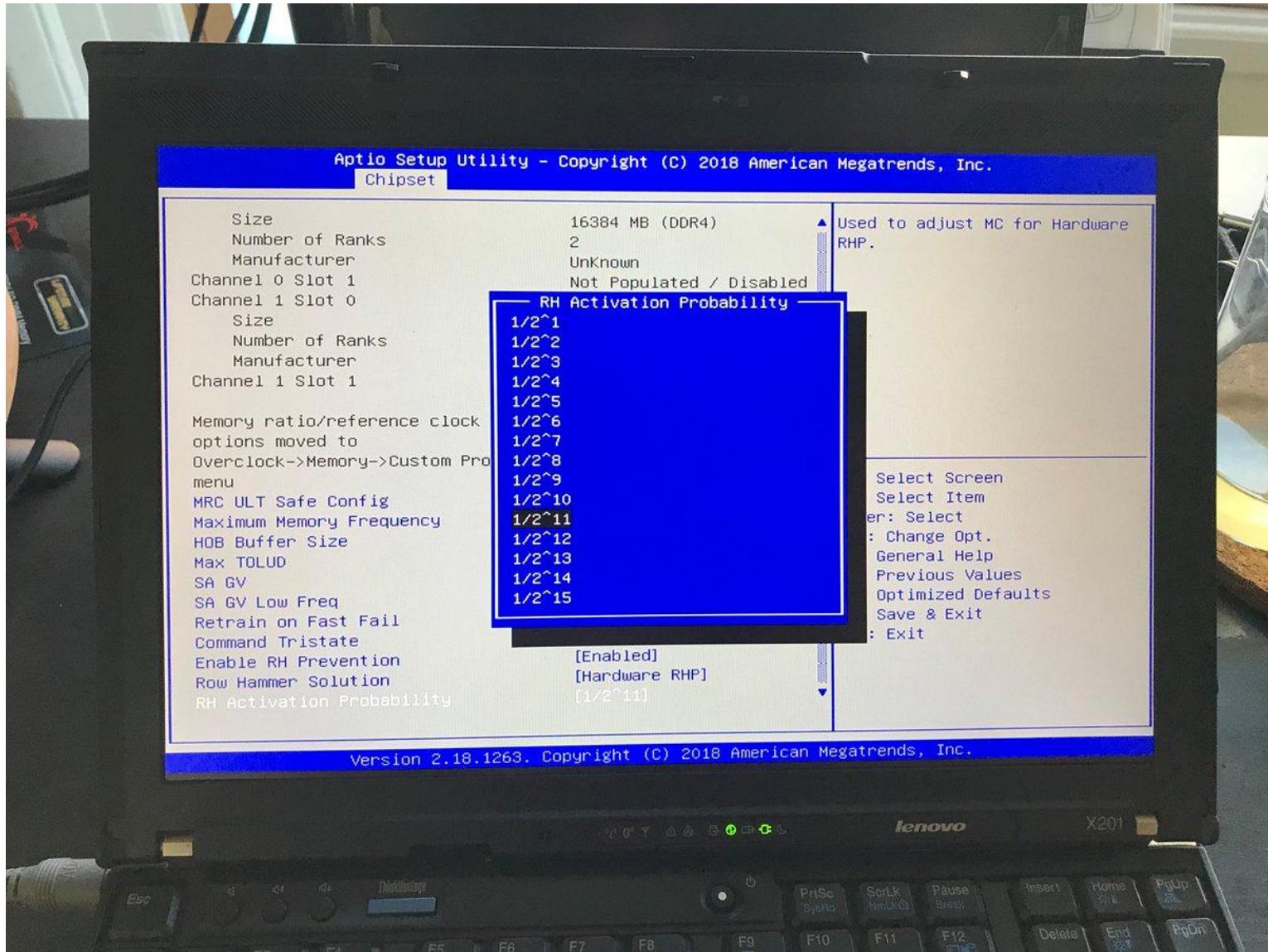
# Requirements for PARA

- If implemented in DRAM chip (done today)
  - Enough slack in timing and refresh parameters
  - Plenty of slack today:
    - Lee et al., “Adaptive-Latency DRAM: Optimizing DRAM Timing for the Common Case,” HPCA 2015.
    - Chang et al., “Understanding Latency Variation in Modern DRAM Chips,” SIGMETRICS 2016.
    - Lee et al., “Design-Induced Latency Variation in Modern DRAM Chips,” SIGMETRICS 2017.
    - Chang et al., “Understanding Reduced-Voltage Operation in Modern DRAM Devices,” SIGMETRICS 2017.
    - Ghose et al., “What Your DRAM Power Models Are Not Telling You: Lessons from a Detailed Experimental Study,” SIGMETRICS 2018.
    - Kim et al., “Solar-DRAM: Reducing DRAM Access Latency by Exploiting the Variation in Local Bitlines,” ICCD 2018.
- If implemented in memory controller
  - Need coordination between controller and DRAM
  - Memory controller should know which rows are physically adjacent

# Probabilistic Activation in Real Life (I)



# Probabilistic Activation in Real Life (II)



# Row Migration-Based RowHammer Defenses

**Key Idea:** Dynamically remap an aggressor row address to a different physical row before a RowHammer bitflip occurs

- Does **not** require refreshing victim rows
- Relocates the aggressor row's data

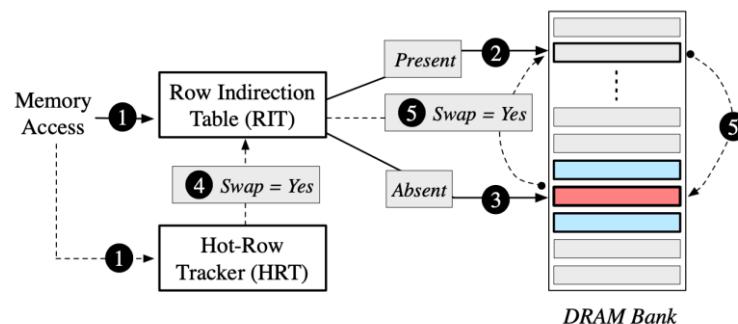


Figure 2: Overview of the Randomized Row Swap (RRS). The Row Indirection Table (RIT) is checked to determine if the access should go to original or remapped location. The Hot-Row Tracker (HRT) identifies rows that must undergo swap.

Saileshwar+, “Randomized Row Swap: Mitigating Row Hammer by Breaking Spatial Correlation between Aggressor and Victim Rows,” ASPLOS’22.

Hassan+, “CROW: A Low-Cost Substrate for Improving DRAM Performance, Energy Efficiency, and Reliability,” ISCA’19.

# Row Migration-Based RowHammer Defenses

## ISCA 2019

### CROW: A Low-Cost Substrate for Improving DRAM Performance, Energy Efficiency, and Reliability

Hasan Hassan<sup>†</sup> Minesh Patel<sup>†</sup> Jeremie S. Kim<sup>†§</sup> A. Giray Yaglikci<sup>†</sup>

Nandita Vijaykumar<sup>†§</sup> Nika Mansouri Ghiasi<sup>†</sup> Saugata Ghose<sup>§</sup> Onur Mutlu<sup>†§</sup>



### Randomized Row-Swap: Mitigating Row Hammer by Breaking Spatial Correlation between Aggressor and Victim Rows

Gururaj Saileshwar\*

Bolin Wang

Moinuddin Qureshi

Prashant J. Nair



### AQUA: Scalable Rowhammer Mitigation by Quarantining Aggressor Rows at Runtime

Anish Saxena

Gururaj Saileshwar

Prashant J. Nair

Moinuddin Qureshi

## HPCA 2023

The 29th IEEE International Symposium on High-Performance Computer Architecture  
(HPCA-29)

**Scalable and Secure Row-Swap:  
Efficient and Safe Row Hammer  
Mitigation in Memory Systems**  
Jeonghyun Woo (University of  
British Columbia),  
Gururaj Saileshwar (Georgia  
Institute of Technology),  
Prashant J. Nair (University of  
British Columbia)

**SHADOW: Preventing Row  
Hammer in DRAM with Intra-  
Subarray Row Shuffling**  
Minbok Wi (Seoul National  
University),  
Jaehyun Park (Seoul National  
University),  
Seoyoung Ko (Seoul National  
University), Michael Jaemin Kim  
(Seoul National University),  
Nam Sung Kim (UIUC),  
Eojin Lee (Inha University),  
Jung Ho Ahn (Seoul National  
University)

# RowHammer in 2023: SK Hynix

## ISSCC 2023 / SESSION 28 / HIGH-DENSITY MEMORIES

### 28.8 A 1.1V 16Gb DDR5 DRAM with Probabilistic-Aggressor Tracking, Refresh-Management Functionality, Per-Row Hammer Tracking, a Multi-Step Precharge, and Core-Bias Modulation for Security and Reliability Enhancement

Woongrae Kim, Chulmoon Jung, Seongnyuh Yoo, Duckhwa Hong,  
Jeongjin Hwang, Jungmin Yoon, Ohyong Jung, Joonwoo Choi, Sanga Hyun,  
Mankeun Kang, Sangho Lee, Dohong Kim, Sanghyun Ku, Donhyun Choi,  
Nogeun Joo, Sangwoo Yoon, Junseok Noh, Byeongyong Go, Cheolhoe Kim,  
Sunil Hwang, Mihyun Hwang, Seol-Min Yi, Hyungmin Kim, Sanghyuk Heo,  
Yeonsu Jang, Kyoungchul Jang, Shinho Chu, Yoonna Oh, Kwidong Kim,  
Junghyun Kim, Soohwan Kim, Jeongtae Hwang, Sangil Park, Junphyo Lee,  
Inchul Jeong, Joohwan Cho, Jonghwan Kim

SK hynix Semiconductor, Icheon, Korea

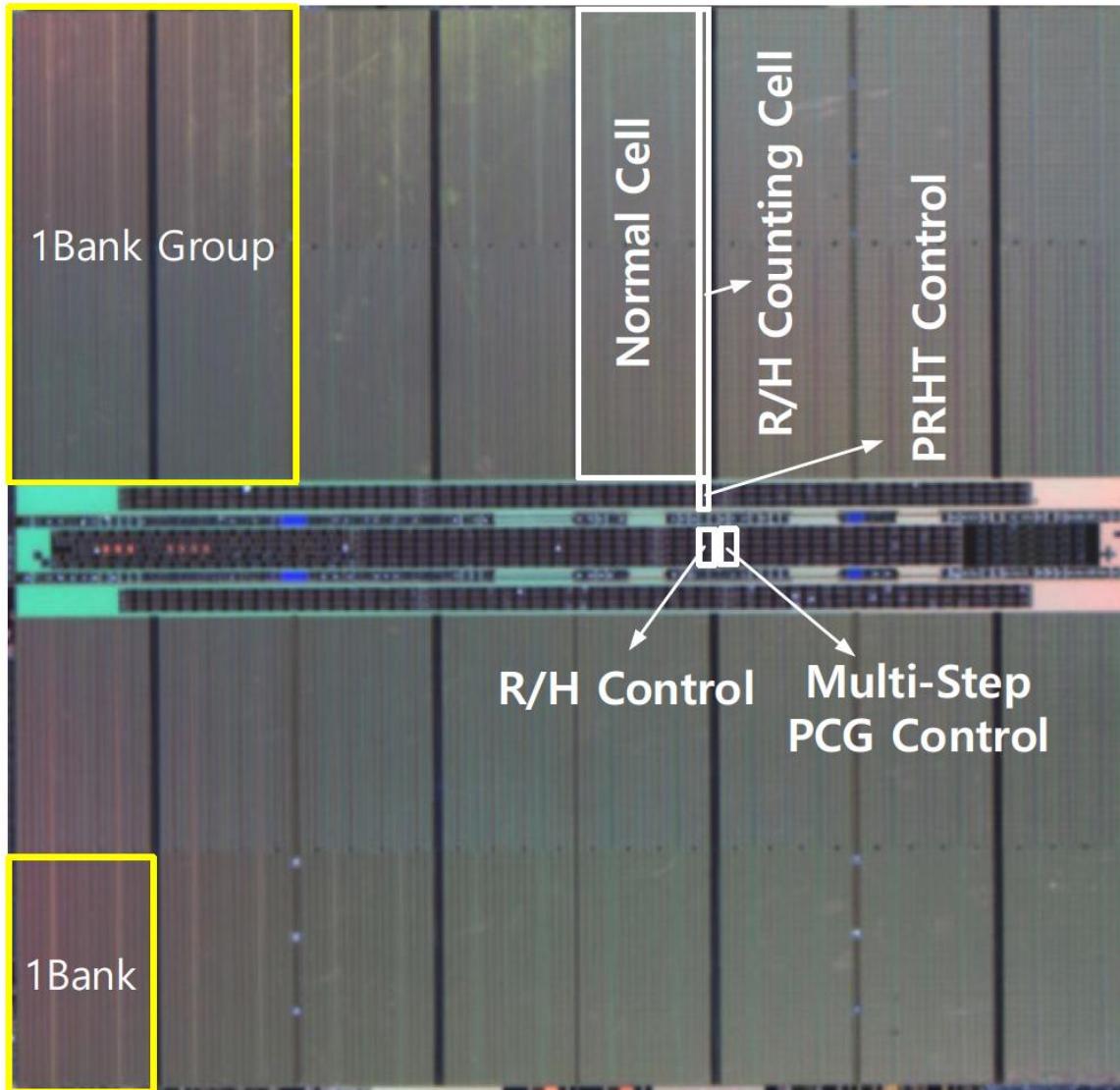


# Industry's RowHammer Solutions (I)

SK hynix Semiconductor, Icheon, Korea

DRAM products have been recently adopted in a wide range of high-performance computing applications: such as in cloud computing, in big data systems, and IoT devices. This demand creates larger memory capacity requirements, thereby requiring aggressive DRAM technology node scaling to reduce the cost per bit [1,2]. However, DRAM manufacturers are facing technology scaling challenges due to row hammer and refresh retention time beyond 1a-nm [2]. Row hammer is a failure mechanism, where repeatedly activating a DRAM row disturbs data in adjacent rows. Scaling down severely threatens reliability since a reduction of DRAM cell size leads to a reduction in the intrinsic row hammer tolerance [2,3]. To improve row hammer tolerance, there is a need to probabilistically activate adjacent rows with carefully sampled active addresses and to improve intrinsic row hammer tolerance [2]. In this paper, row-hammer-protection and refresh-management schemes are presented to guarantee DRAM security and reliability despite the aggressive scaling from 1a-nm to sub 10-nm nodes. The probabilistic-aggressor-tracking scheme with a refresh-management function (RFM) and per-row hammer tracking (PRHT) improve DRAM resilience. A multi-step precharge reinforces intrinsic row-hammer tolerance and a core-bias modulation improves retention time: even in the face of cell-transistor degradation due to technology scaling. This comprehensive scheme leads to a reduced probability of failure, due to row hammer attacks, by 93.1% and an improvement in retention time by 17%.

# Industry's RowHammer Solutions (II)



ISSCC 2023 / SESSION 28 / HIGH-DENSITY MEMORIES /

28.8 A 1.1V 16Gb DDR5 DRAM with Probabilistic-Aggressor Tracking, Refresh-Management Functionality, Per-Row Hammer Tracking, a Multi-Step Precharge, and Core-Bias Modulation for Security and Reliability Enhancement

Woongrae Kim, Chulmoon Jung, Seongnyuh Yoo, Duckhwa Hong, Jeongjin Hwang, Jungmin Yoon, Ohyong Jung, Joonwoo Choi, Sanga Hyun, Mankeun Kang, Sangho Lee, Dohong Kim, Sanghyun Ku, Donhyun Choi, Nogeuon Joo, Sangwoo Yoon, Junseok Noh, Byeongyong Go, Cheolhoe Kim, Sunil Hwang, Mihyun Hwang, Seol-Min Yi, Hyungmin Kim, Sanghyuk Heo, Yeonsu Jang, Kyoungchul Jang, Shinho Chu, Yoonna Oh, Kwidong Kim, Junghyun Kim, Soohwan Kim, Jeongtae Hwang, Sangil Park, Junphyo Lee, Inchul Jeong, Joohwan Cho, Jonghwan Kim

SK hynix Semiconductor, Icheon, Korea

# DRAM Controller

(本節內容改自 Prof. Onur Mutlu, “Computer Architecture,” 12<sup>th</sup> Lecture, Fall 2023 課程講義)

# DRAM Types

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

# DRAM Types (circa 2015)

<i>Segment</i>	<i>DRAM Standards &amp; Architectures</i>
Commodity	DDR3 (2007) [14]; DDR4 (2012) [18]
Low-Power	LPDDR3 (2012) [17]; LPDDR4 (2014) [20]
Graphics	GDDR5 (2009) [15]
Performance	eDRAM [28], [32]; RLDRAM3 (2011) [29]
3D-Stacked	WIO (2011) [16]; WIO2 (2014) [21]; MCDRAM (2015) [13]; HBM (2013) [19]; HMC1.0 (2013) [10]; HMC1.1 (2014) [11]
Academic	SBA/SSA (2010) [38]; Staged Reads (2012) [8]; RAIDR (2012) [27]; SALP (2012) [24]; TL-DRAM (2013) [26]; RowClone (2013) [37]; Half-DRAM (2014) [39]; Row-Buffer Decoupling (2014) [33]; SARP (2014) [6]; AL-DRAM (2015) [25]

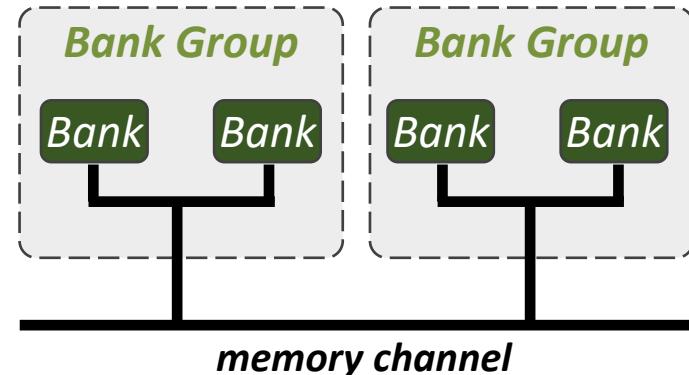
Table 1. Landscape of DRAM-based memory

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

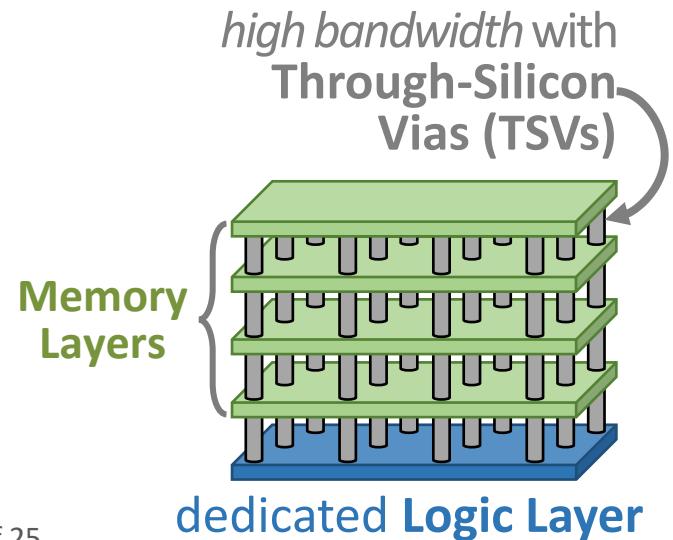
# Modern DRAM Types: Comparison to

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

- Bank groups

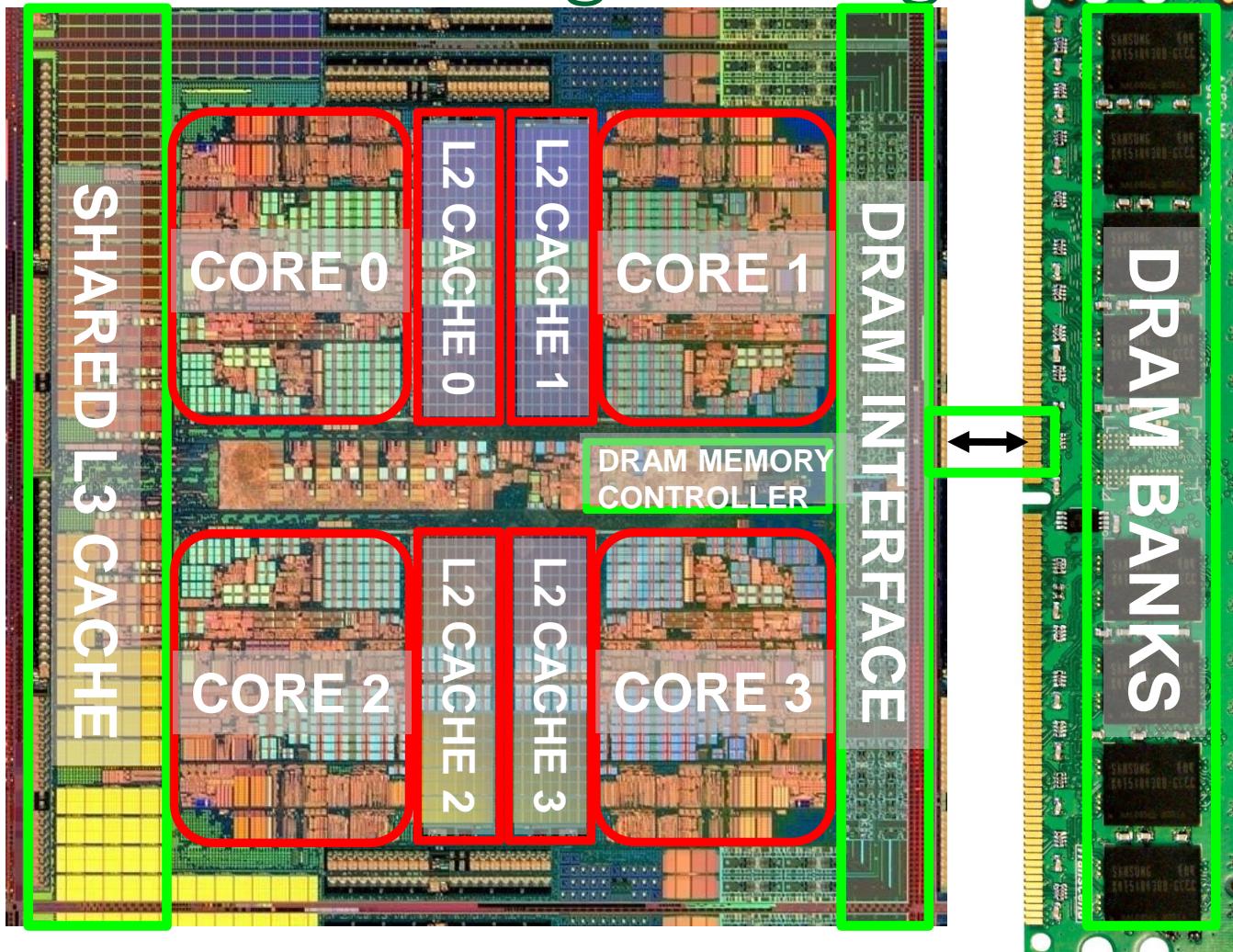


- 3D-stacked DRAM



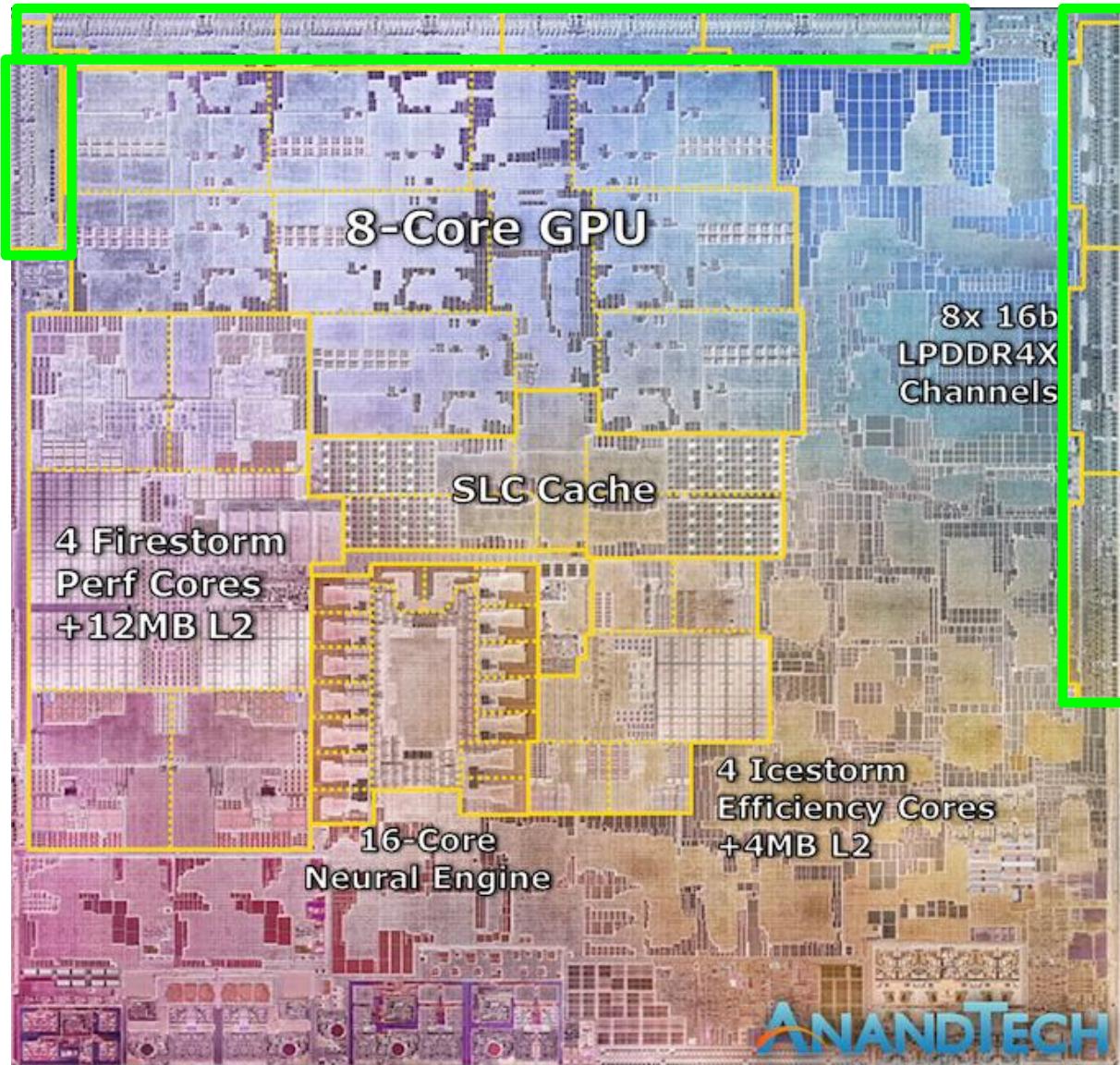
# DRAM Control Logic Is Large

Multi-Core  
Chip



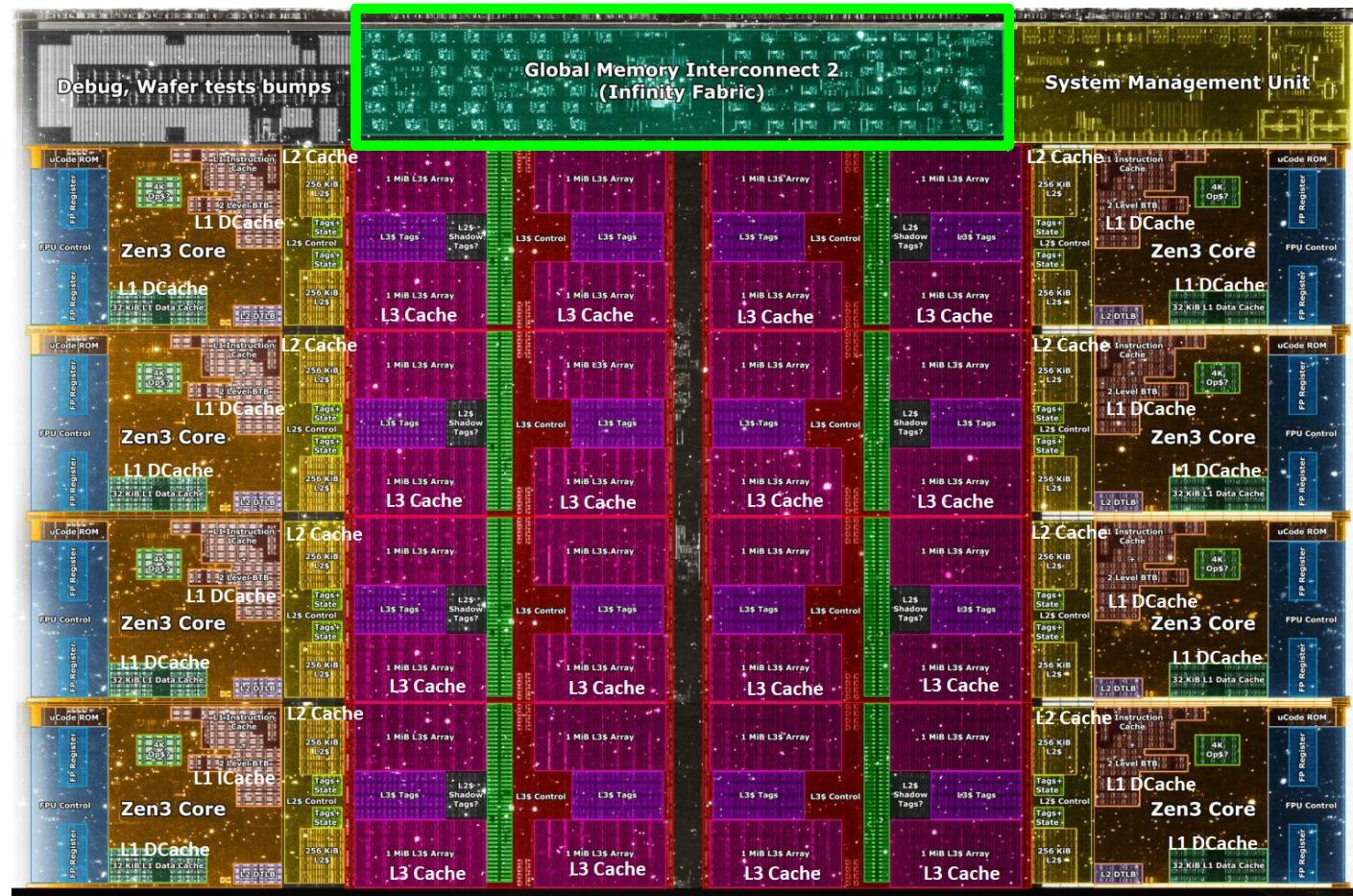
\*Die photo credit: AMD Barcelona

# DRAM Control Logic Is Large



Apple M1,  
2021

# DRAM Control Logic Is Large



Core Count:  
8 cores/16 threads

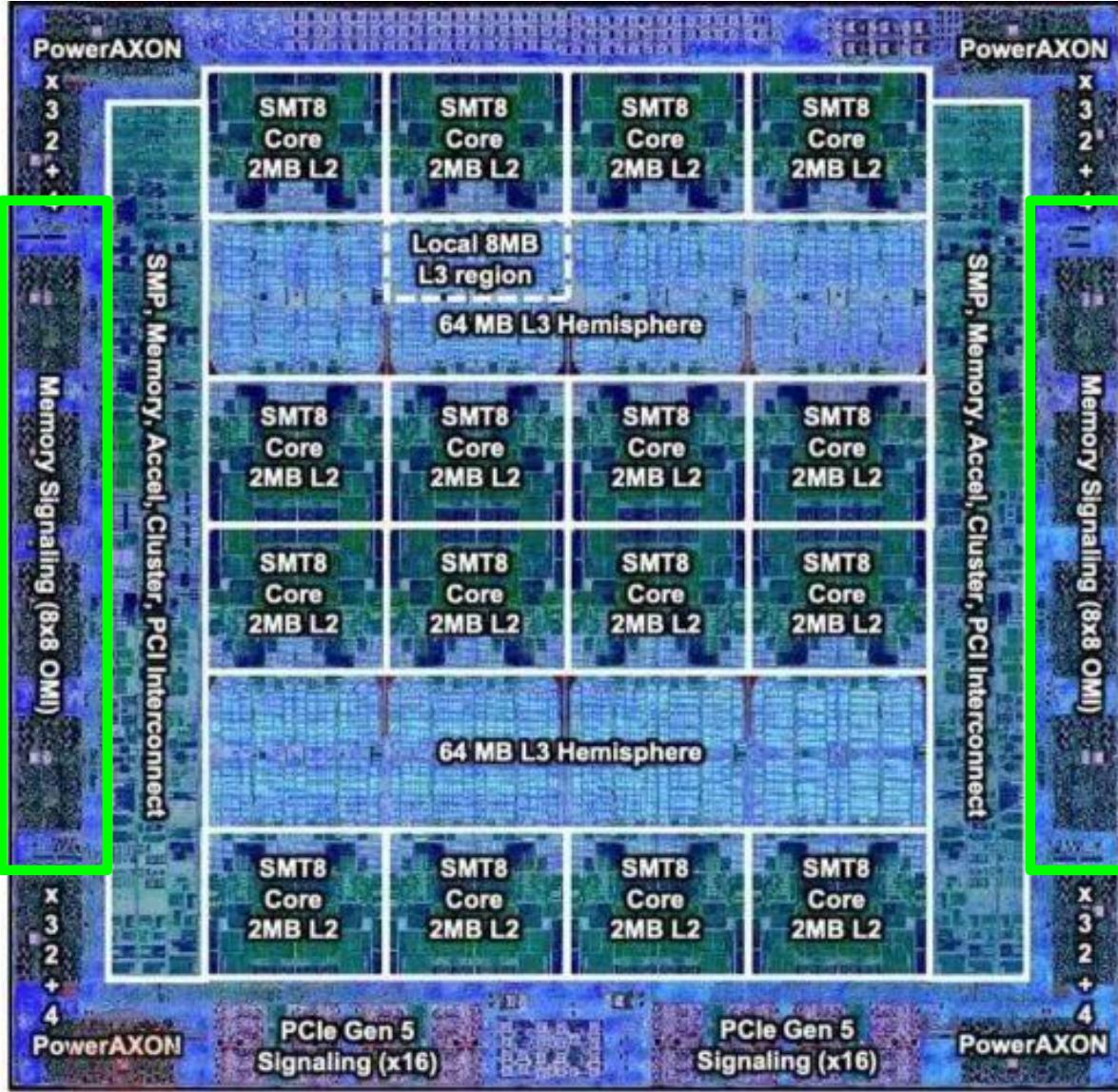
L1 Caches:  
32 KB per core

L2 Caches:  
512 KB per core

L3 Cache:  
32 MB shared

AMD Ryzen 5000, 2020

# DRAM Control Logic Is Large



IBM POWER10,  
2020

Cores:  
15-16 cores,  
8 threads/core

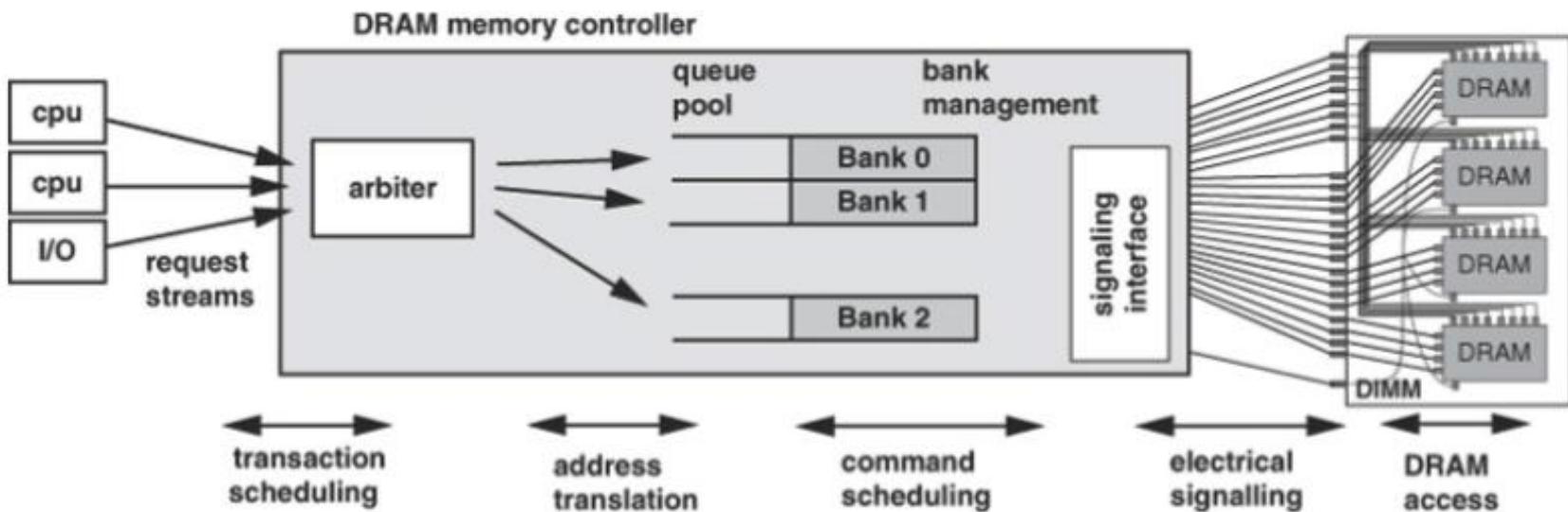
L2 Caches:  
2 MB per core

L3 Cache:  
120 MB shared

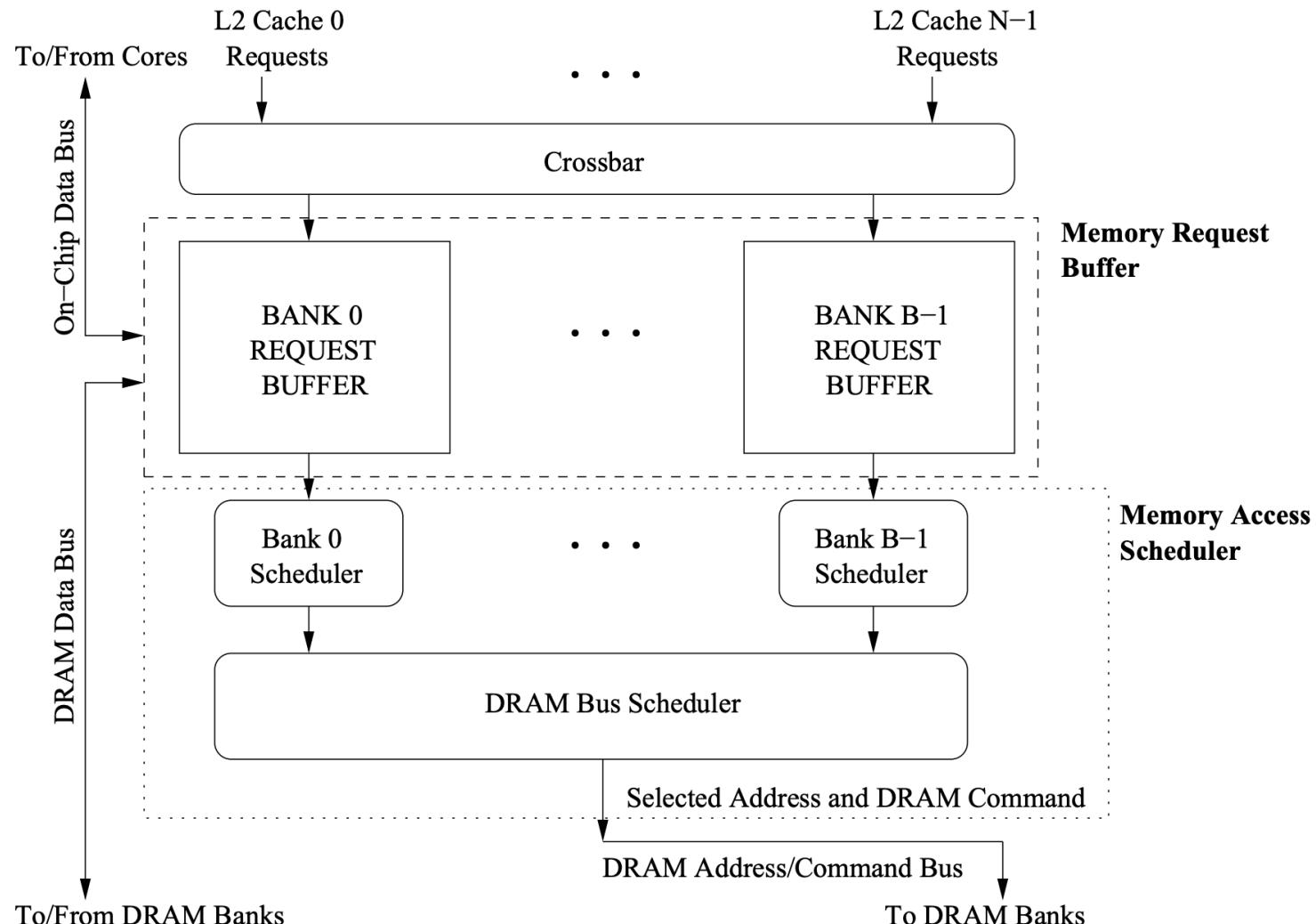
# DRAM Controller: Functions

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

# A Modern DRAM Controller (I)



# A Modern DRAM Controller (II)



# DRAM Scheduling Policies (I)

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

Goal: Maximize row buffer hit rate → **maximize DRAM throughput**

  - Actually, scheduling is done at the **command level**
    - Column commands (read/write) prioritized over row commands (activate/precharge)
    - Within each group (e.g., bank), older commands prioritized over younger ones

# DRAM Bank Operation

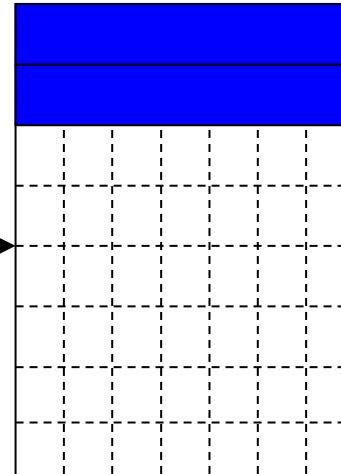
Access Address:

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

Row address 0



Columns

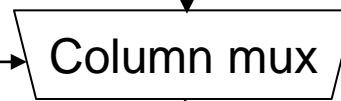


Rows

Row 1

Row Buffer CONFLICT !

Column address 05



Data

# DRAM Scheduling Policies (II)

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

# Row Buffer Management Policies

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

# Open vs. Closed Row Policies

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

# DRAM Power Management

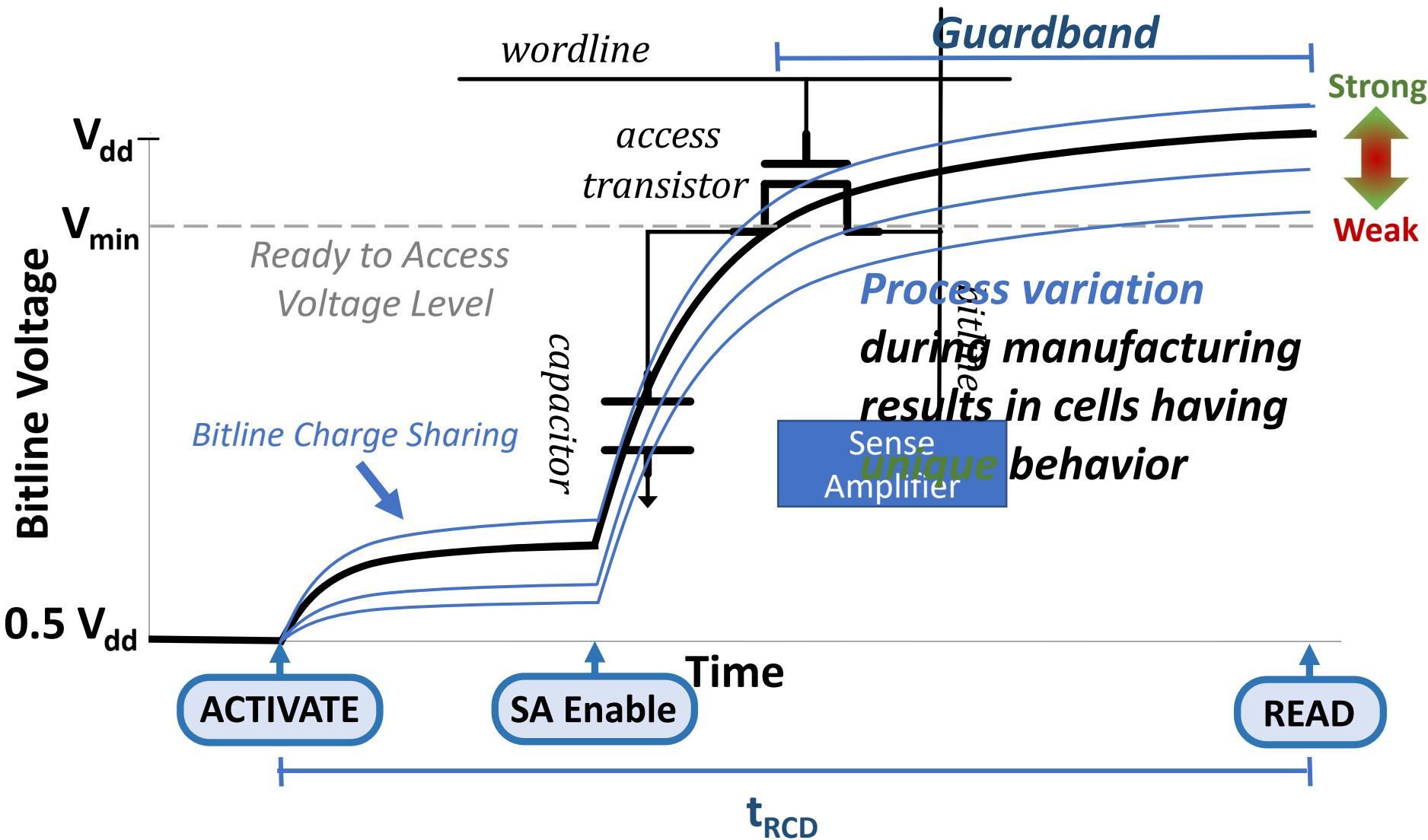
- DRAM chips have power modes
- Idea: When not accessing a chip power it down
- Power states
  - Active (highest power)
  - All banks idle
  - Power-down
  - Self-refresh (lowest power)
- Tradeoff: State transitions incur latency during which the chip cannot be accessed

# Difficulty of DRAM Control

# Why Are DRAM Controllers Difficult to Design?

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

# Why Timing Constraints?



# Many DRAM Timing Constraints

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

Table 4. DDR3 1600 DRAM timing specifications

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

# More on DRAM Operation

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

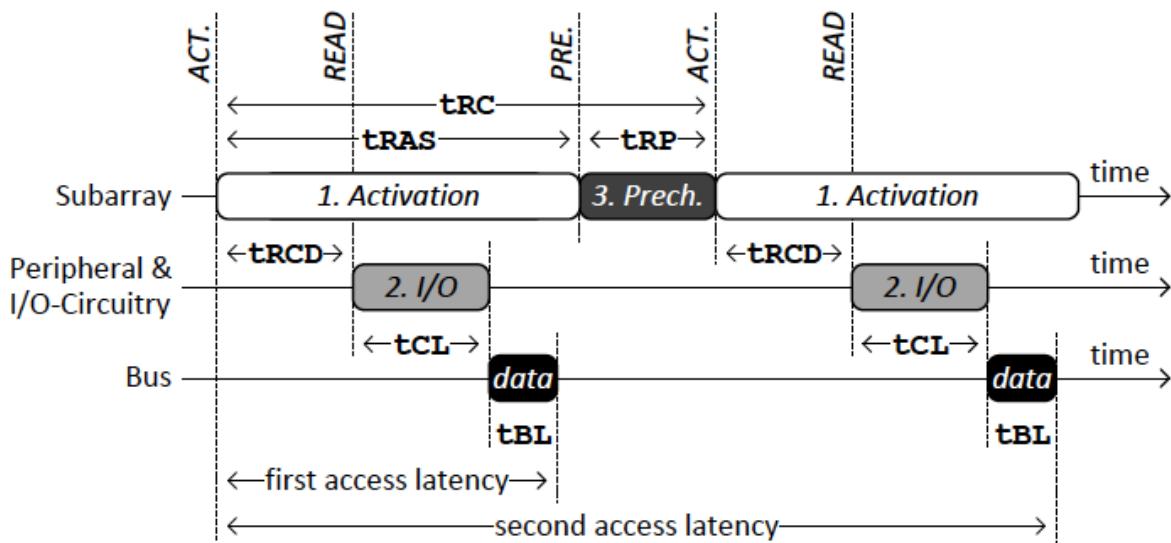


Figure 5. Three Phases of DRAM Access

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

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

# Why So Many Timing Constraints? (I)

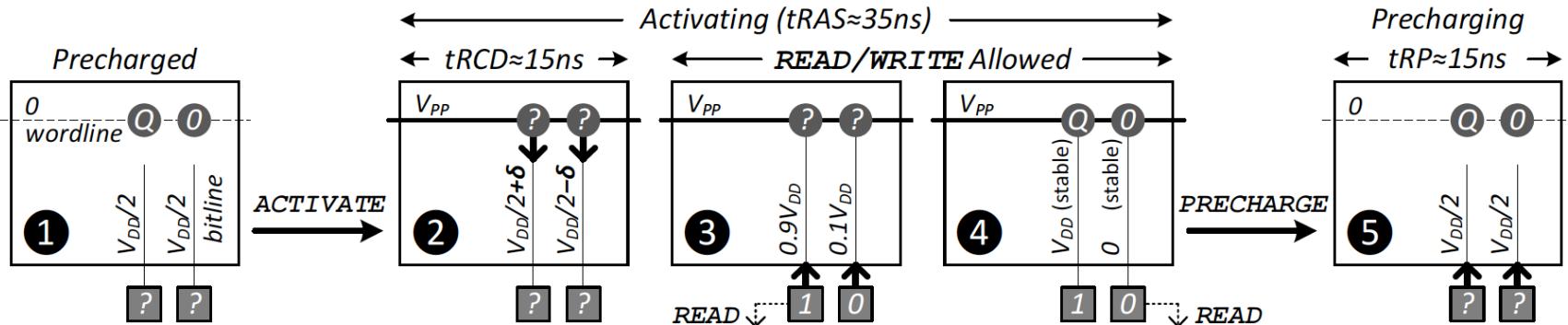


Figure 4. DRAM bank operation: Steps involved in serving a memory request [17] ( $V_{PP} > V_{DD}$ )

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

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

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

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

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

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

# Why So Many Timing Constraints? (II)

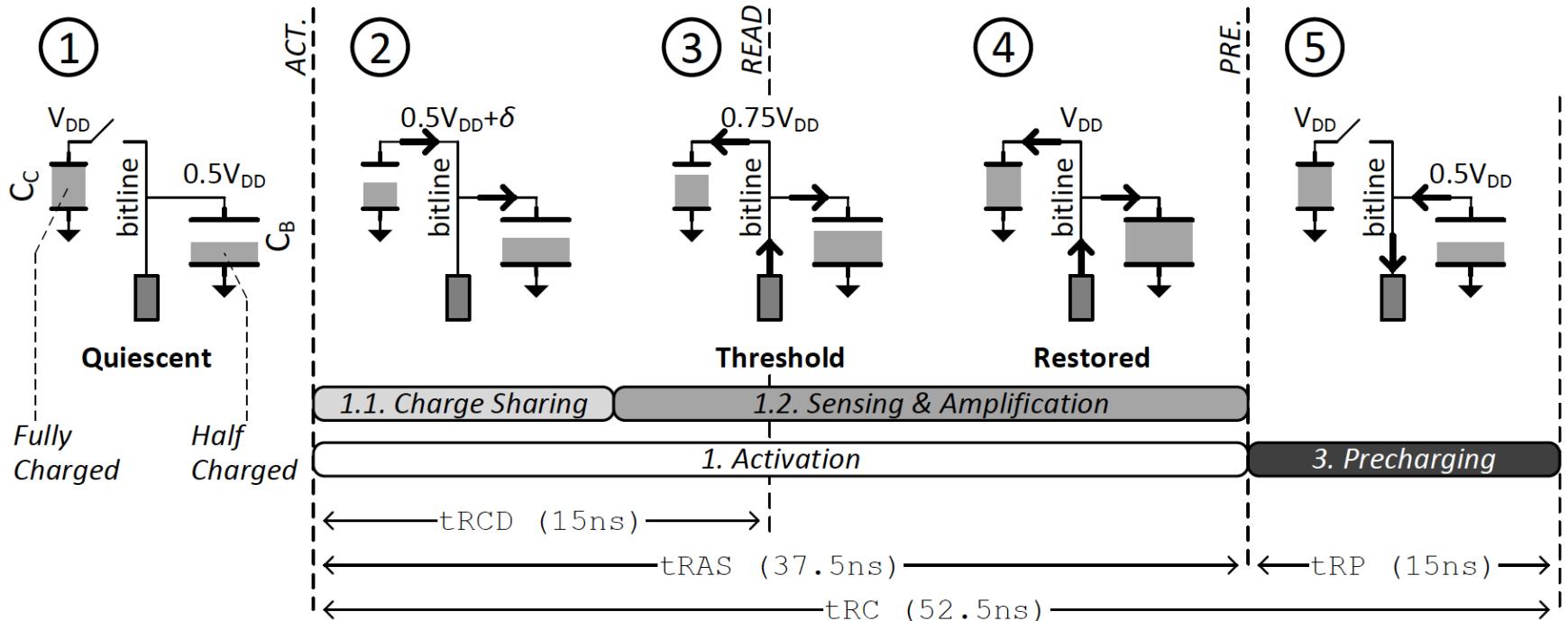


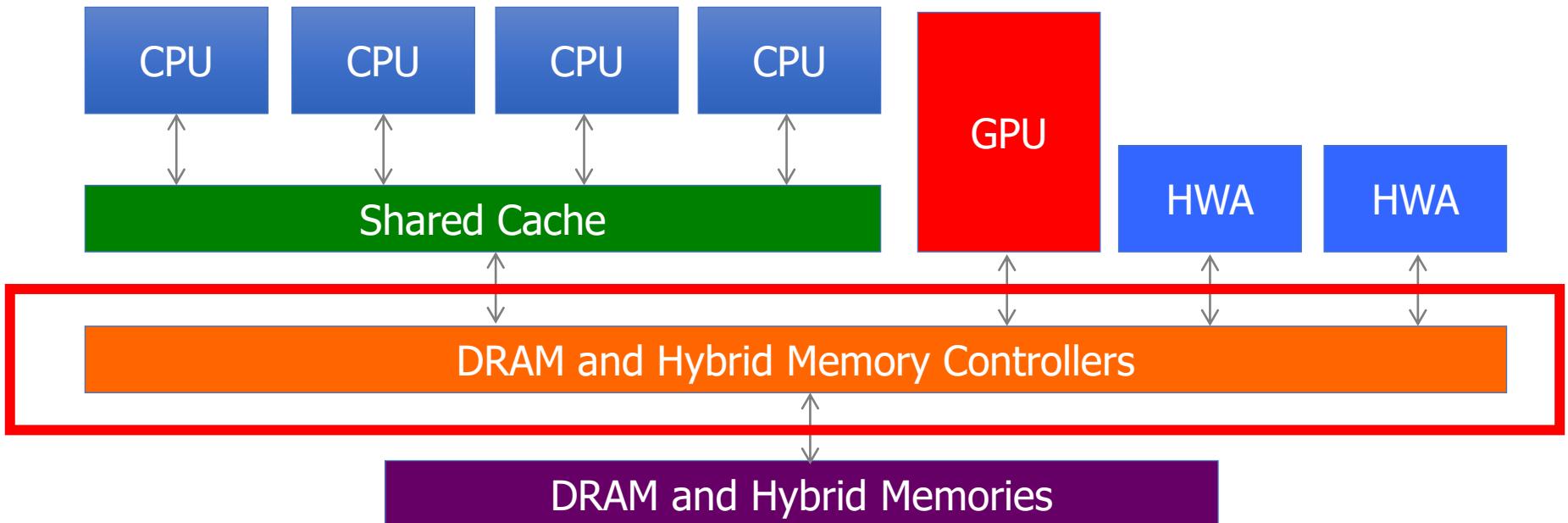
Figure 6. Charge Flow Between the Cell Capacitor ( $C_C$ ), Bitline Parasitic Capacitor ( $C_B$ ), and the Sense-Amplifier ( $C_B \approx 3.5C_C$  [39])

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

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

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

# DRAM Controller Design Is Becoming More Difficult



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