

N26F300

VLSI SYSTEM DESIGN

(GRADUATE LEVEL)

Memory Subsystems: Dynamic Random Memories
(DRAMs)

Outline

2

- DRAM Overview
- DRAM Device and Structure
- DRAM Subsystem Organization

Reference: “Memory Systems Cache, DRAM, DISK” by Bruce Jacob, Spencer W. Ng and David T. Wang

[part of slides are adopted from MOE materials “DRAM Sub-system,” by Prof. CT Huang, NTHU, 2013]

3

DRAM Overview

Why DRAM matters?

4

- Patented in 1968 by Dennard
- Significantly cheaper than SRAM
 - ▣ Higher density than SRAMs
 - 1T vs. 6T
 - ▣ A bit is represented by a high or low charge on the capacitor
- Significantly slower than SRAM
 - ▣ Off-chip (external) vs. on-chip (embedded)
- Disadvantages
 - ▣ Longer access times
 - ▣ Leaky, needs to be refreshed
 - ▣ Cannot be easily integrated with CMOS

Basics of DRAM

DRAM

- DRAM (Dynamic RAM)
- Used mostly in main mem.
- Capacitor + 1 transistor/bit
- Need refresh every 4-8 ms
 - ▣ 5% of total time
- Read is destructive (need for write-back)
- Access time < cycle time (because of writing back)
- Density (25-50):1 to SRAM
- Address lines multiplexed
 - ▣ pins are scarce!

SRAM

- SRAM (Static RAM)
- Used mostly in caches (L, D, TLB, BTB)
- 1 flip-flop (4-6 transistors) per bit
- Read is not destructive
- Access time = cycle time
- Speed (8-16):1 to DRAM
- Address lines not multiplexed
 - ▣ high speed of decoding imp.

Cost of Different Kinds of Storage Devices

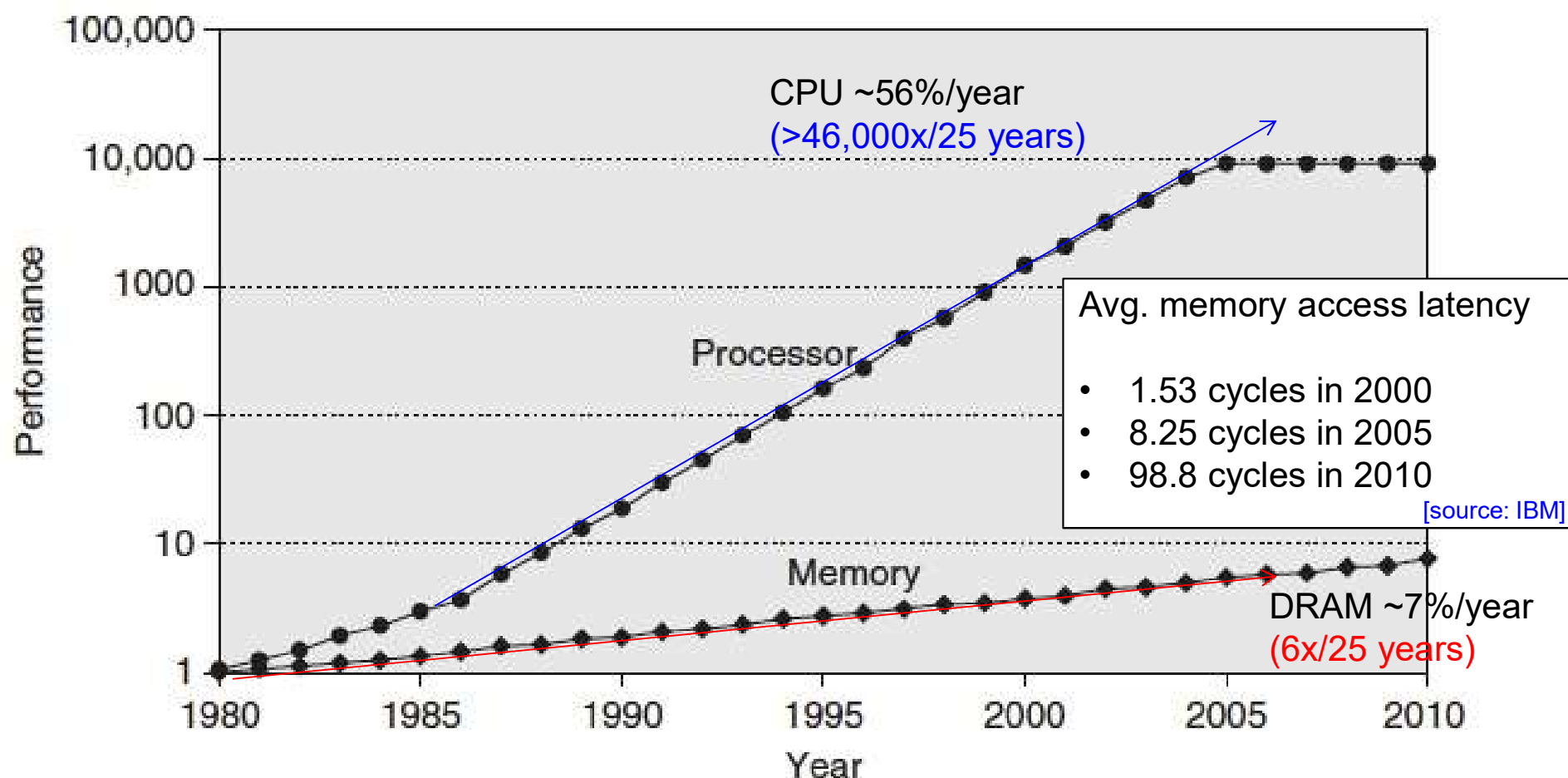
6

Technology	Typical Access	\$ per GB in 2008
SRAM	0.5 – 2.5 ns	\$2000 - \$5000
DRAM	50 – 70 ns	\$20 - \$75
Magnetic Disk	5 – 20 ms	\$0.2 - \$2

Processor-DRAM Memory Gap

7

Processor-DRAM Performance Gap grows ~50% / year



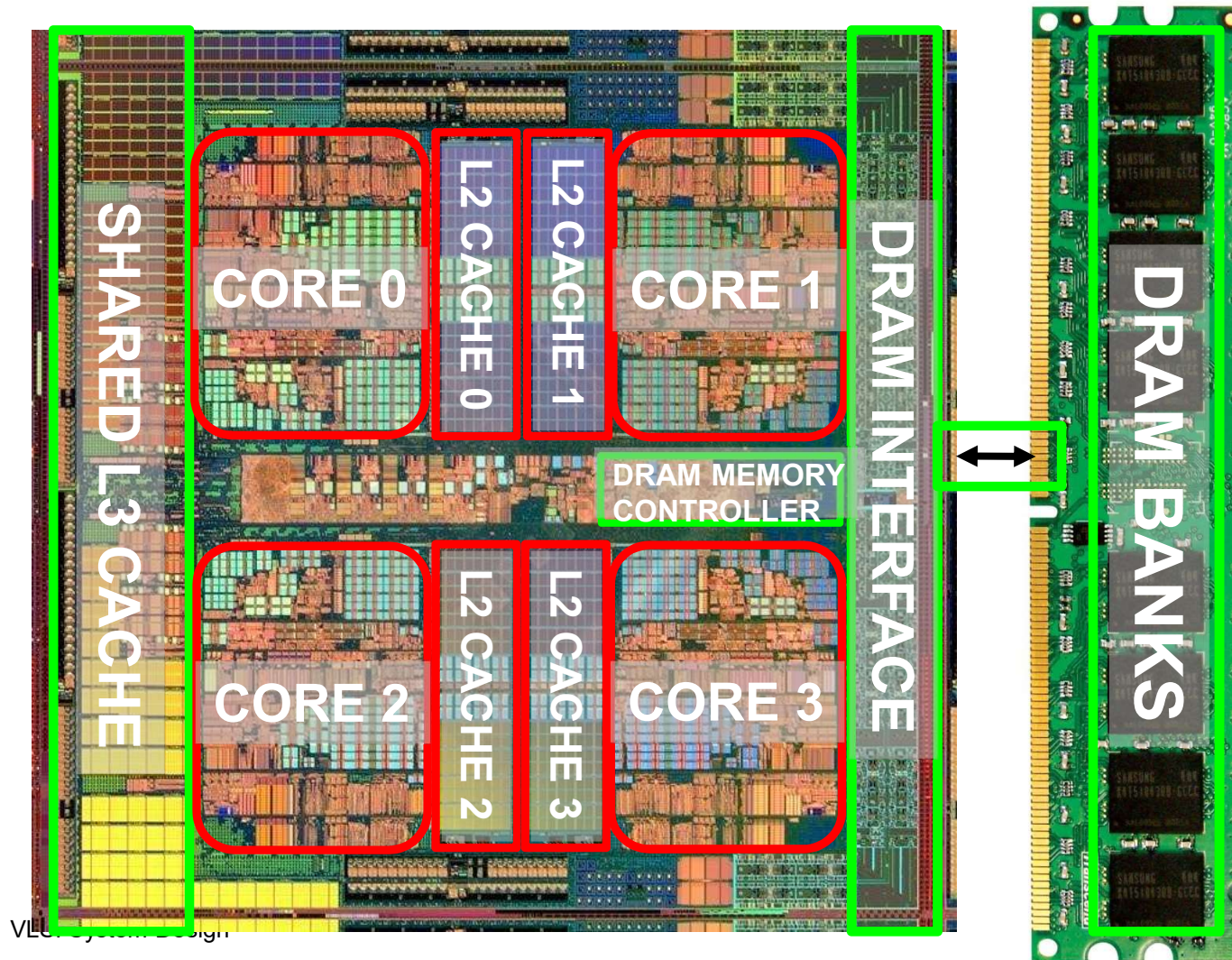
Memory Issues

8

- Latency
 - ▣ Time to move through the longest circuit path (from the start of request to the response)
- Bandwidth
 - ▣ Number of bits transported at one time
- Capacity
 - ▣ Size of memory
- Energy
 - ▣ Cost of accessing memory (to read and write)

Main Memory in The System

9

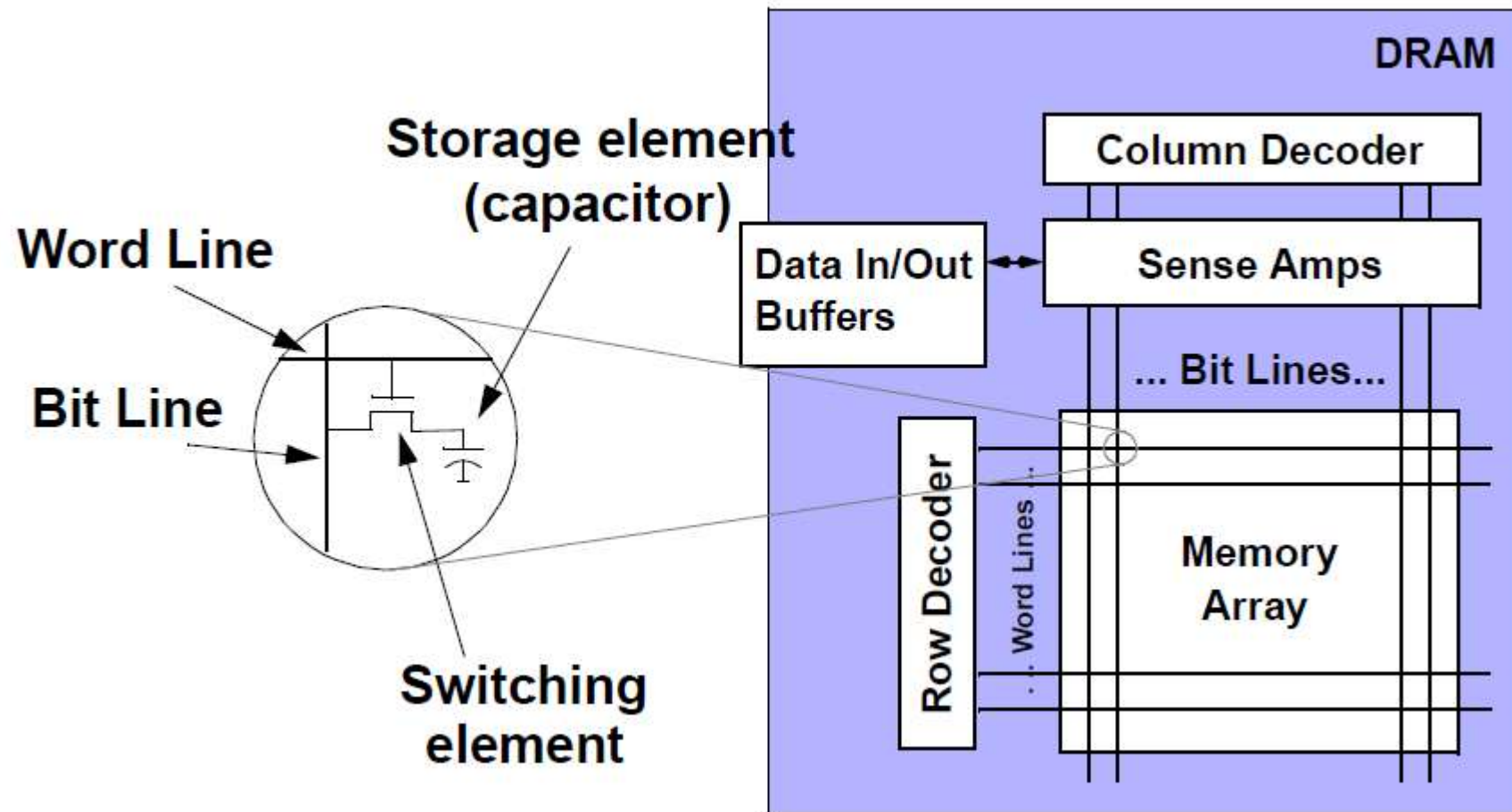


10

DRAM Device & Structure

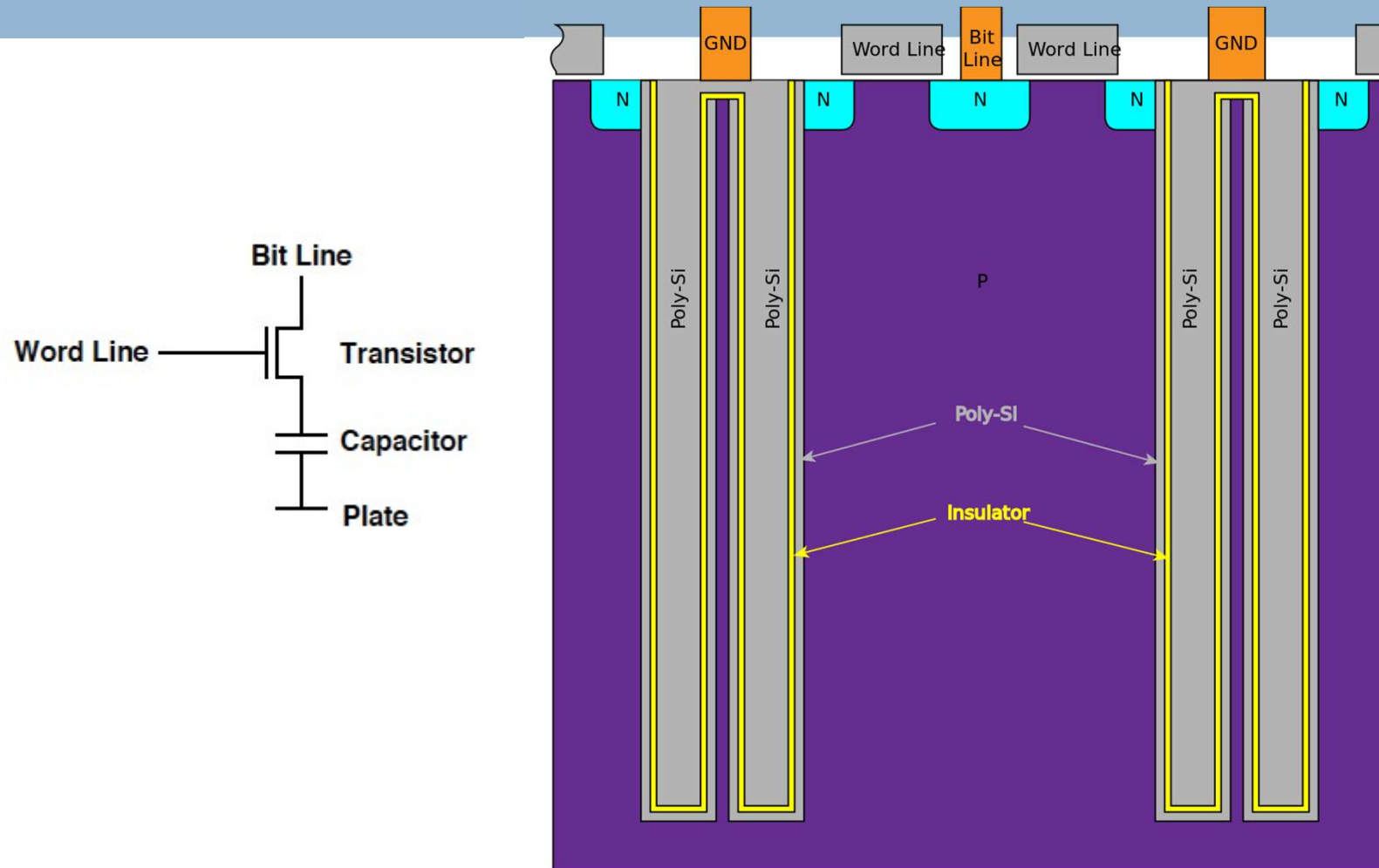
DRAM Cell

11



Trench DRAM Cell

12



1-Transistor Memory Cell (DRAM)

13

Write:

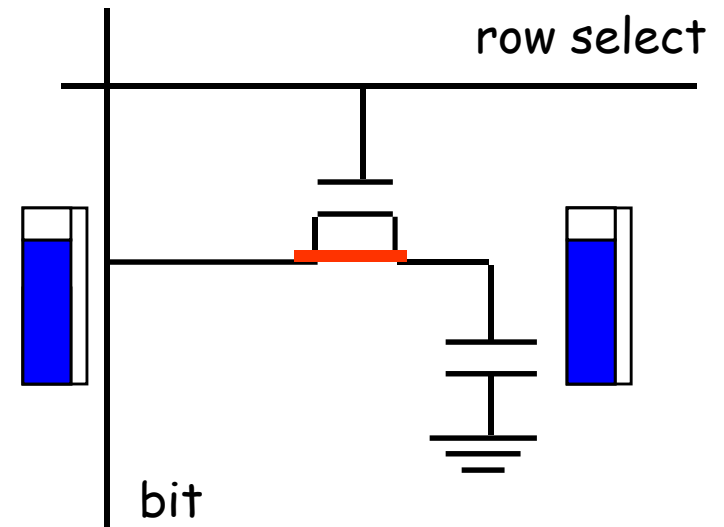
- 1. Drive bit line
- 2. Select row

Read:

- 1. Precharge bit line to $V_{dd}/2$
- 2. Select row
- 3. Cell and bit line share charges
 - Minute voltage changes on the bit line
- 4. Sense (fancy sense amp)
 - Can detect changes of ~ 1 million electrons
- 5. Write: restore the value

Refresh

- 1. Just do a dummy read to every cell.



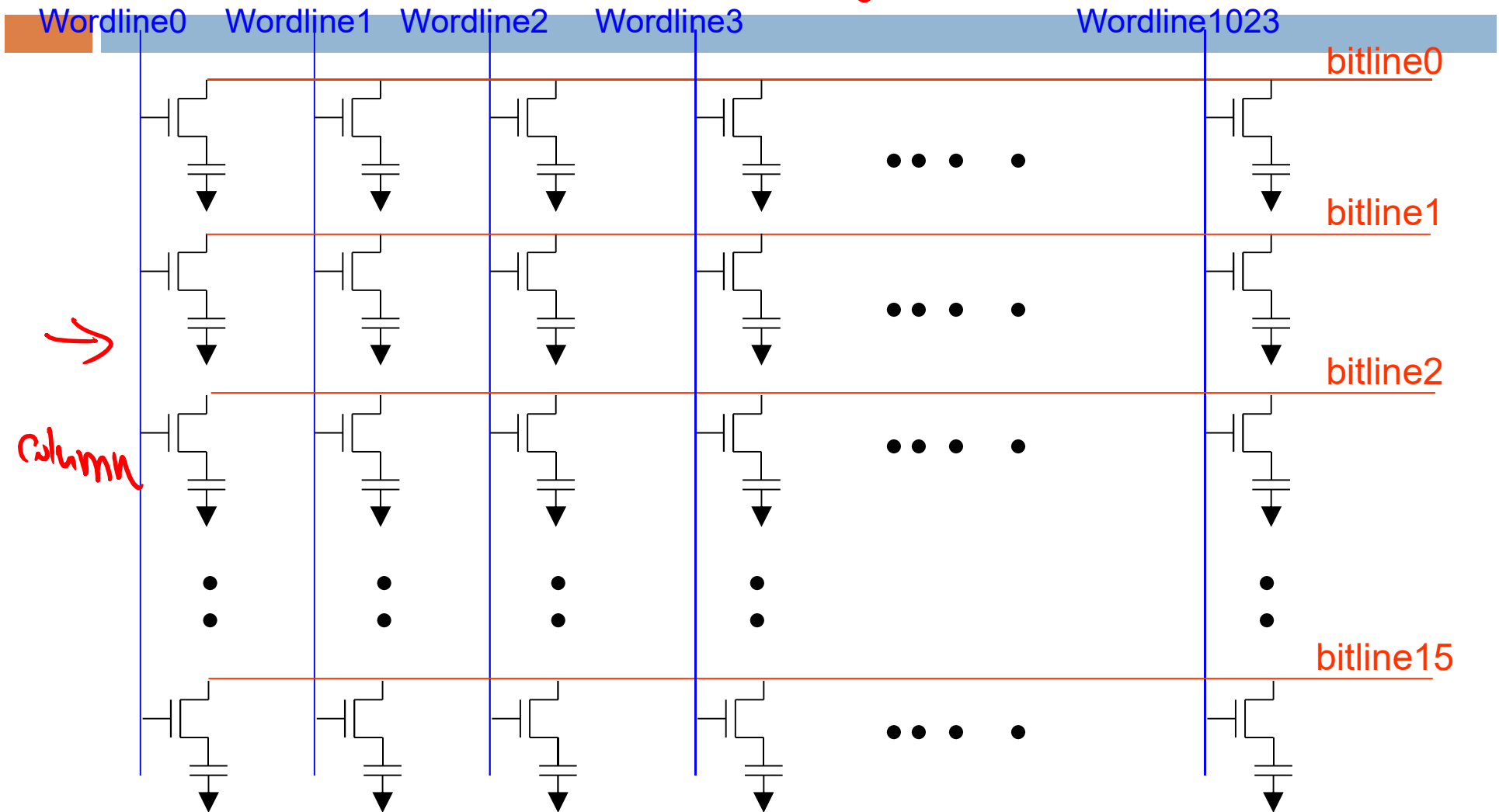
Read is really a read followed by a restoring write

DRAM Cell Properties

14

- Voltage swing on bitline is small
 - ▣ Design target:
 - Bitline capacitance as small as possible; bit cell capacitance as large as possible to increase charge transfer
- Read is destructive
 - ▣ Part of read cycle is used to restore level inside of bit cell capacitor
- Capacitor leaks
 - ▣ Periodical refresh is mandatory
- Noise sources in DRAM are word line to bit line coupling, bit line to bit line coupling

DRAM Cell Array

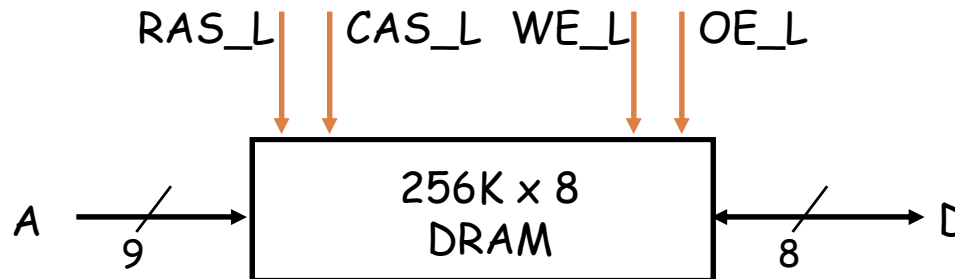


16



Logic Diagram of a Typical DRAM

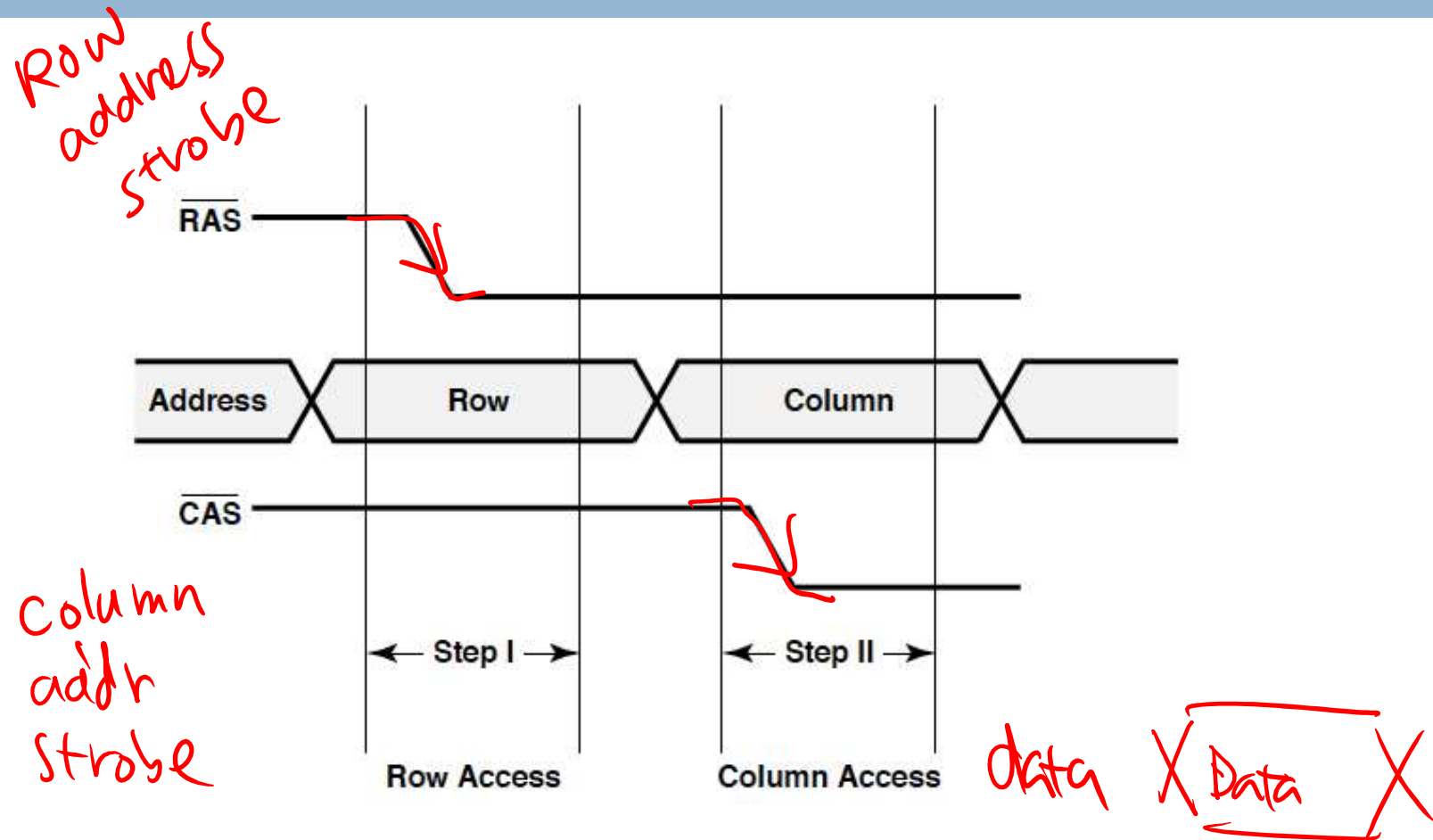
17



- Control Signals (RAS_L, CAS_L, WE_L, OE_L) are all active low
- Din and Dout are combined (D):
 - WE_L is asserted (Low), OE_L is disasserted (High)
 - D serves as the data input pin
 - WE_L is disasserted (High), OE_L is asserted (Low)
 - D is the data output pin
- Row and column addresses share the same pins (A)
 - RAS_L goes low: Pins A are latched in as row address
 - CAS_L goes low: Pins A are latched in as column address
 - RAS/CAS edge-sensitive

Simple DRAM Access Timing

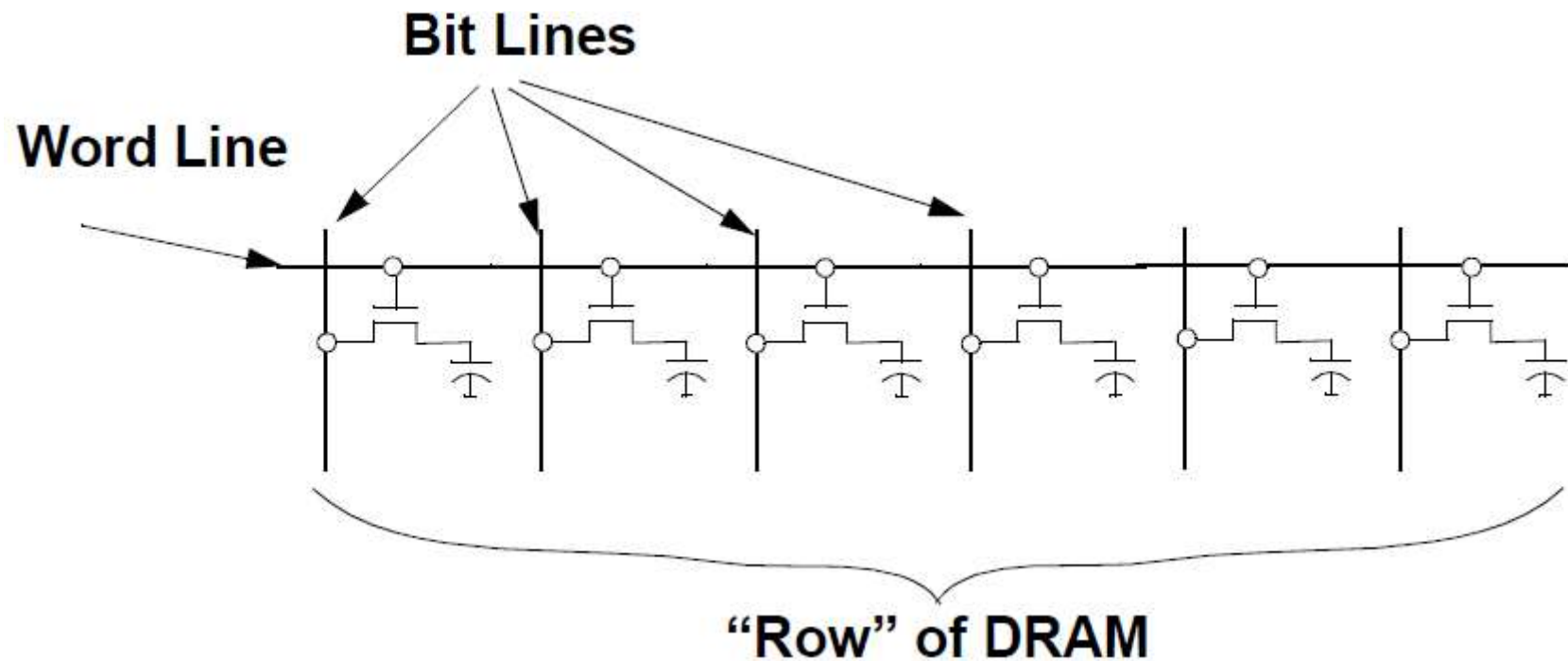
18



Row

19

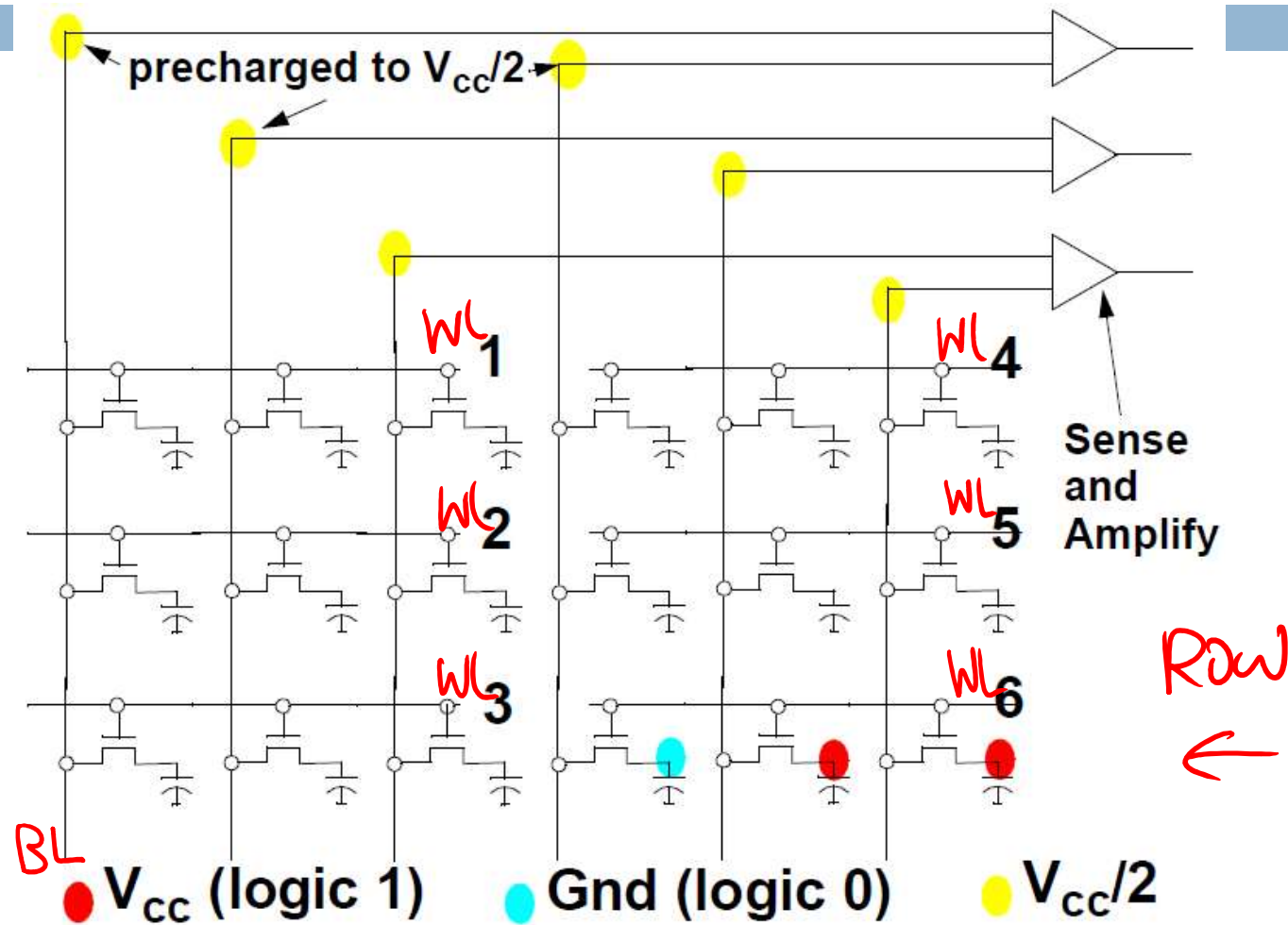
- Row Size: 8 Kb @ 256 Mb SDRAM(synchronous DRAM)



Page size: # of bits per row i.e., # of bits loaded into the sense Amps.

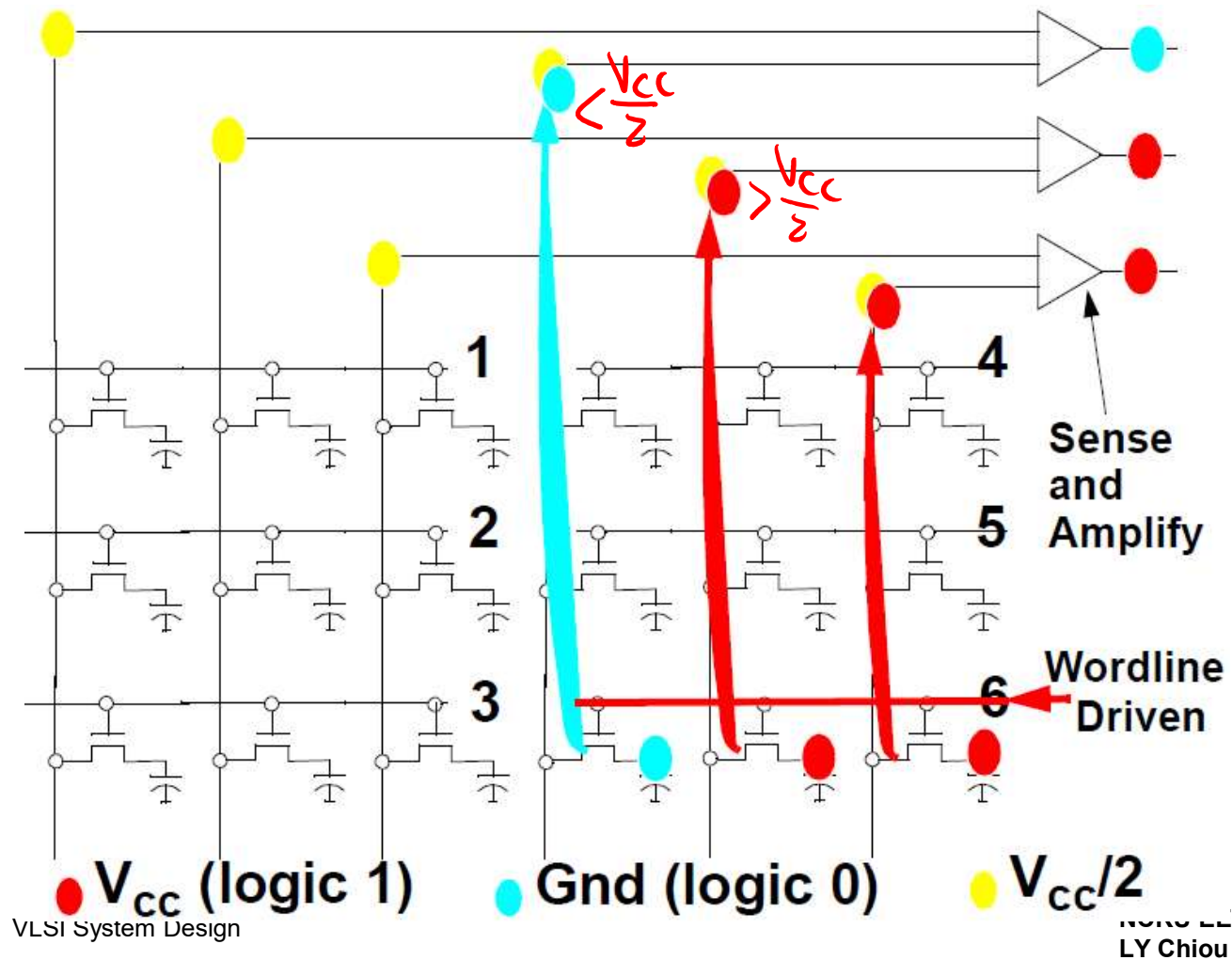
Precharge and Sense Amp

20



Destructive Read

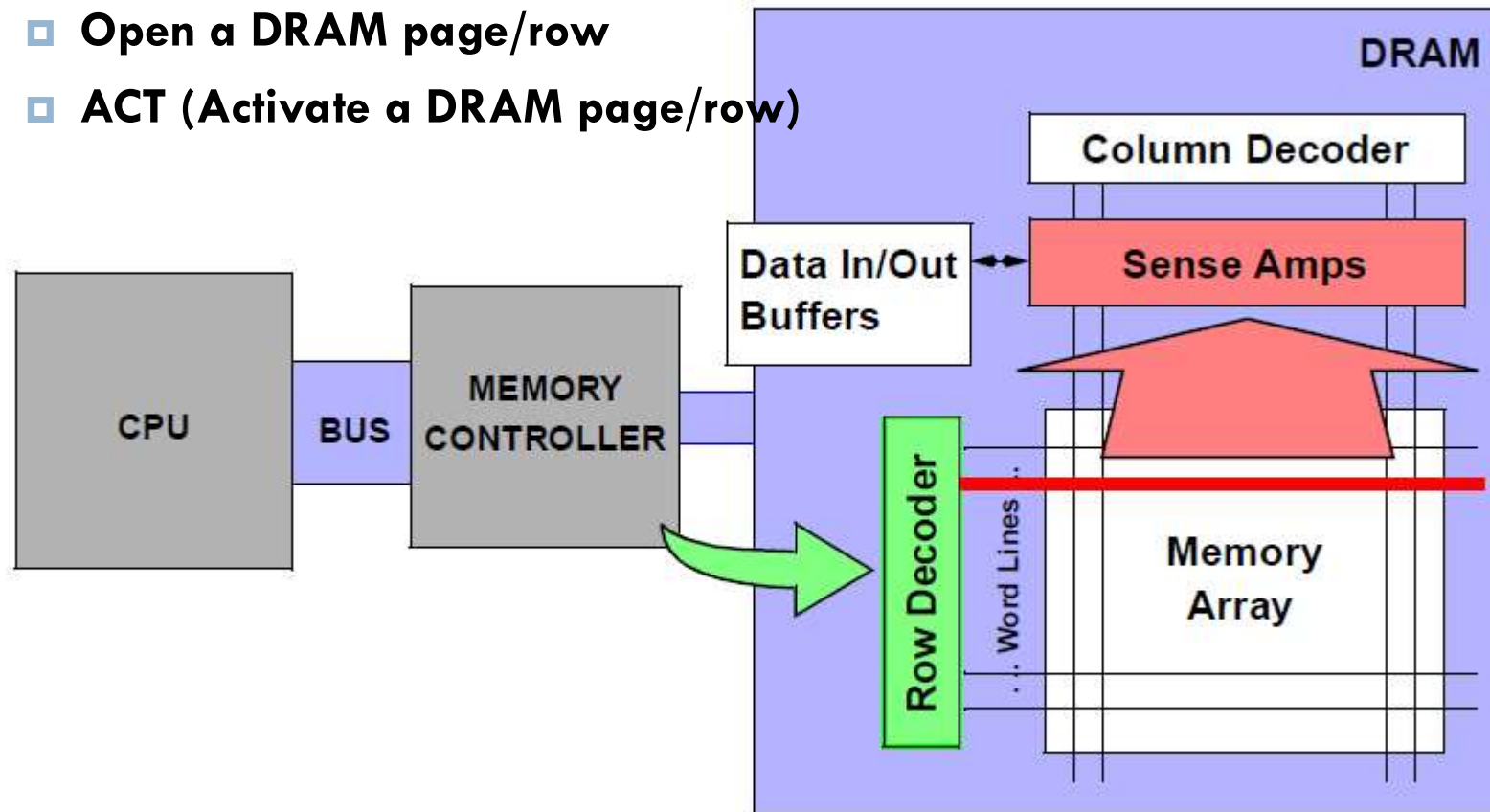
21



Row Access

22

- **RAS (Row Address Strobe)**
 - ▣ Open a DRAM page/row
 - ▣ ACT (Activate a DRAM page/row)

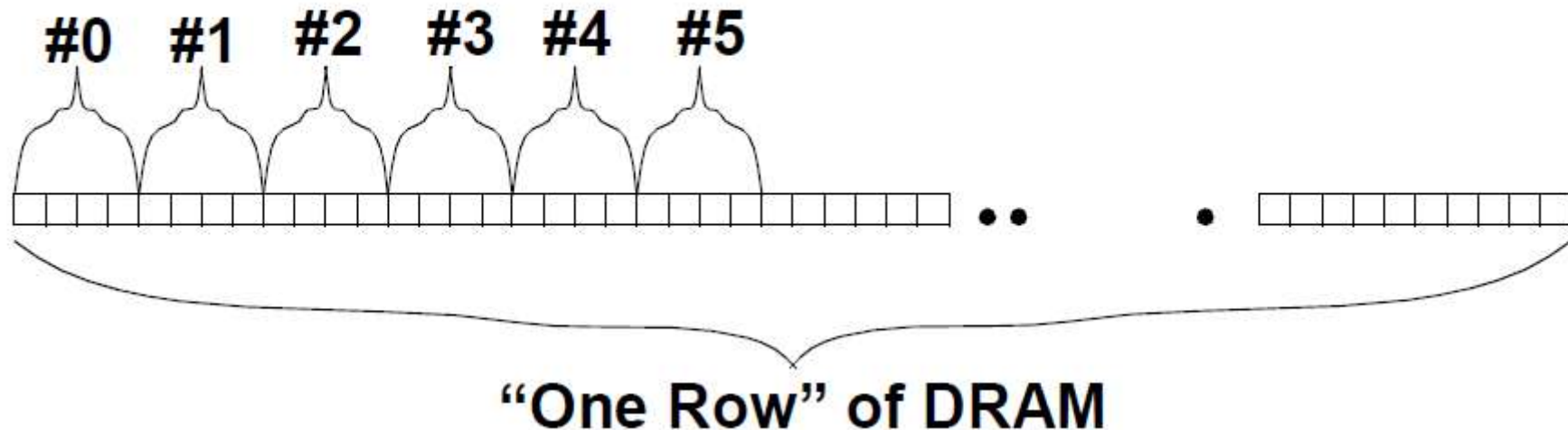


Column

23

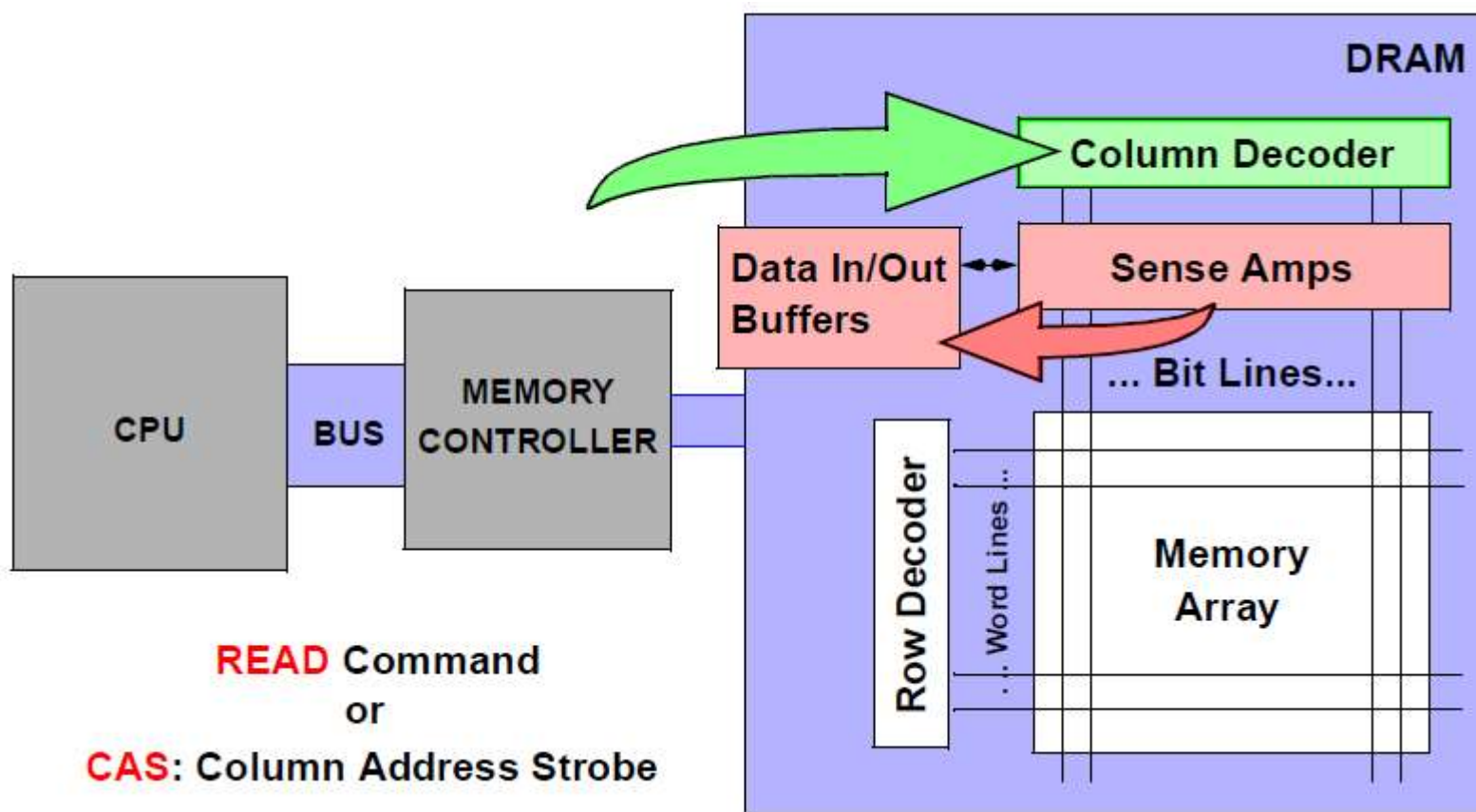
- Smallest addressable quantity of DRAM on chip
- For SDRAM (synchronous DRAM)
 - ▣ column size == chip data bus width (4, 8, 16, 32)
 - ▣ get “n” columns per access. $n = (1, 2, 4, 8)$

4 bit wide columns



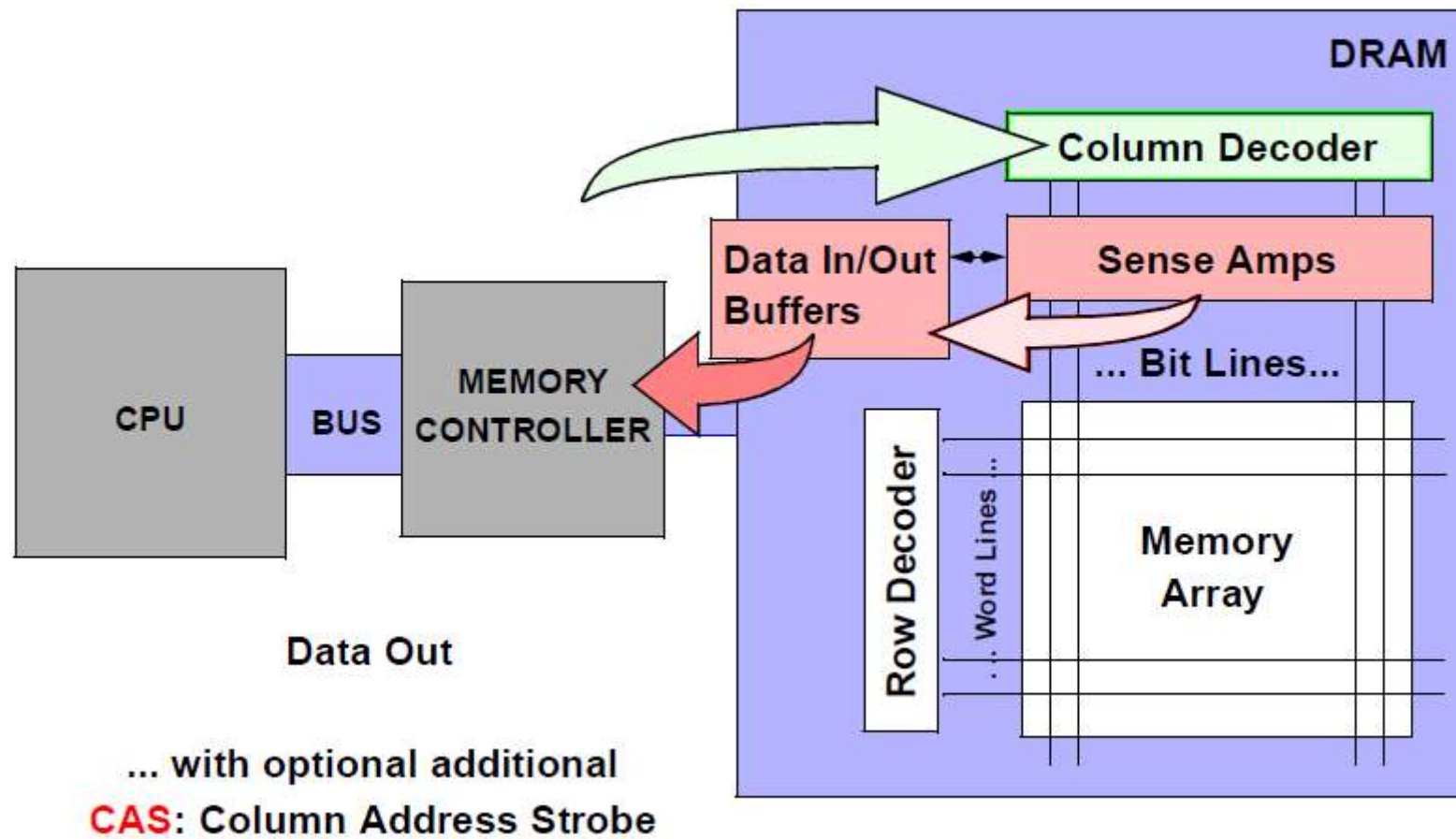
Column Access

24



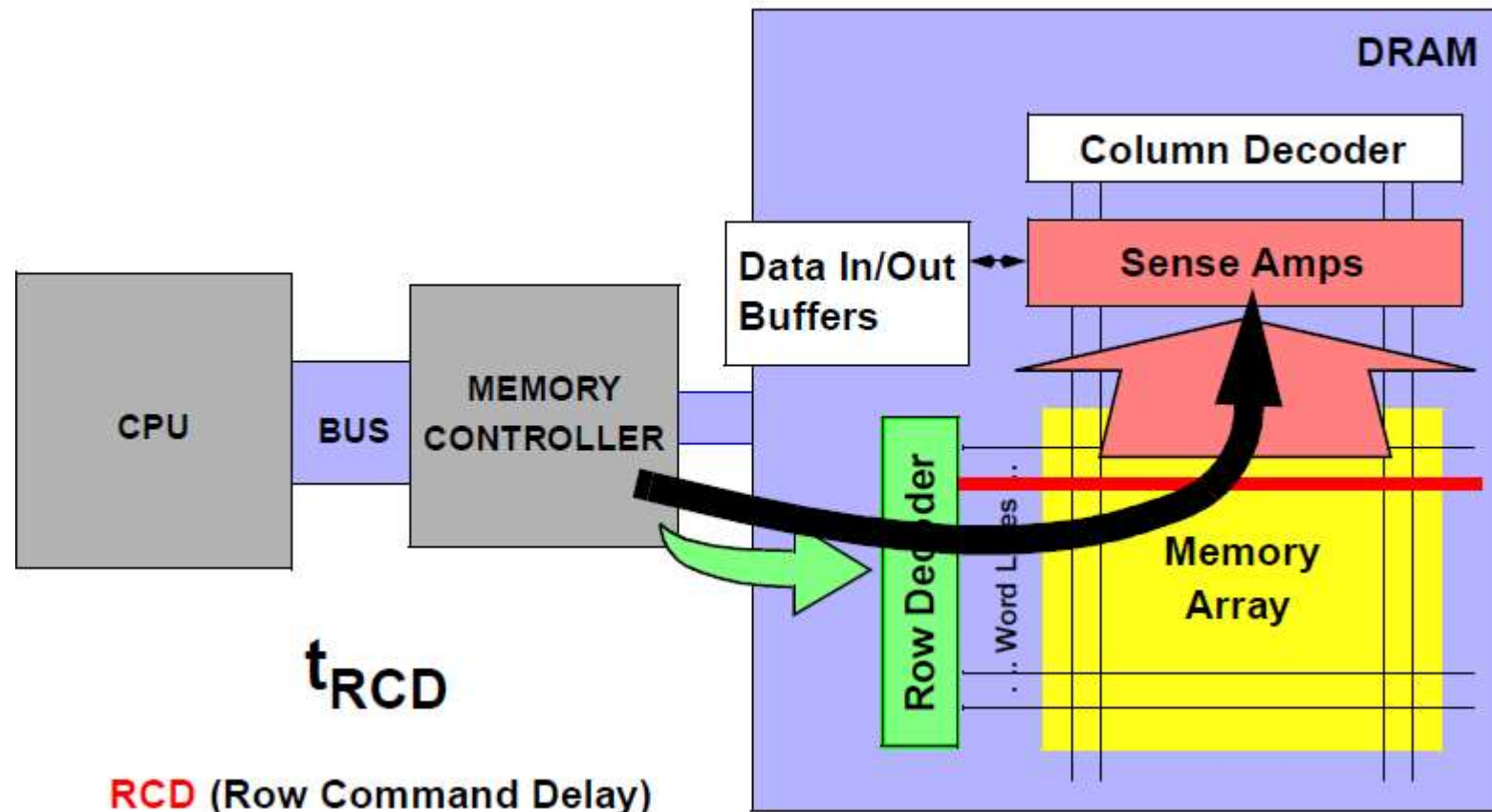
Data Out

25



How fast can I move data from DRAM cell to sense amp?

26



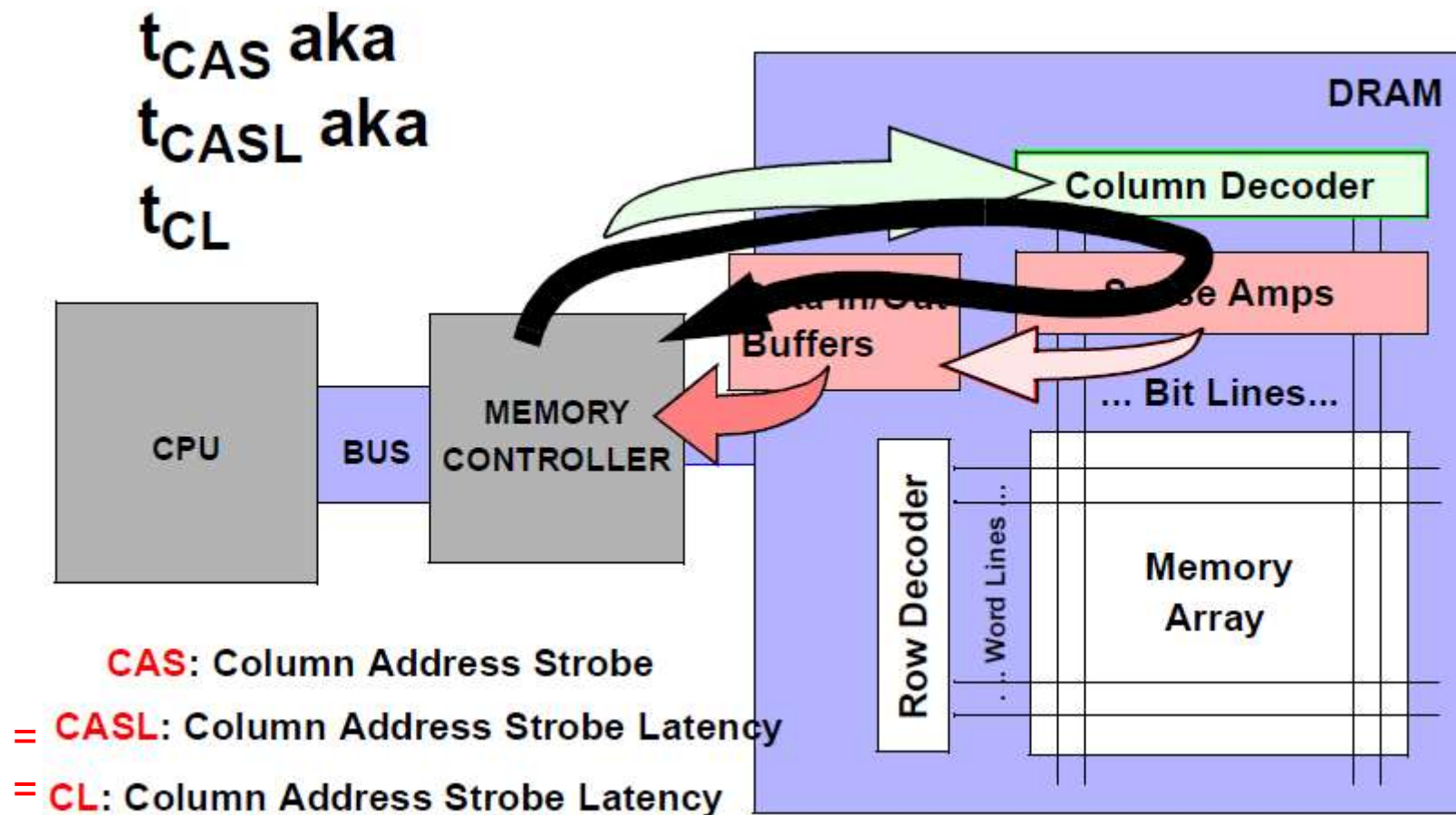
RCD (Row Command Delay)

VLSI System Design

NCKU EE
LY Chiou

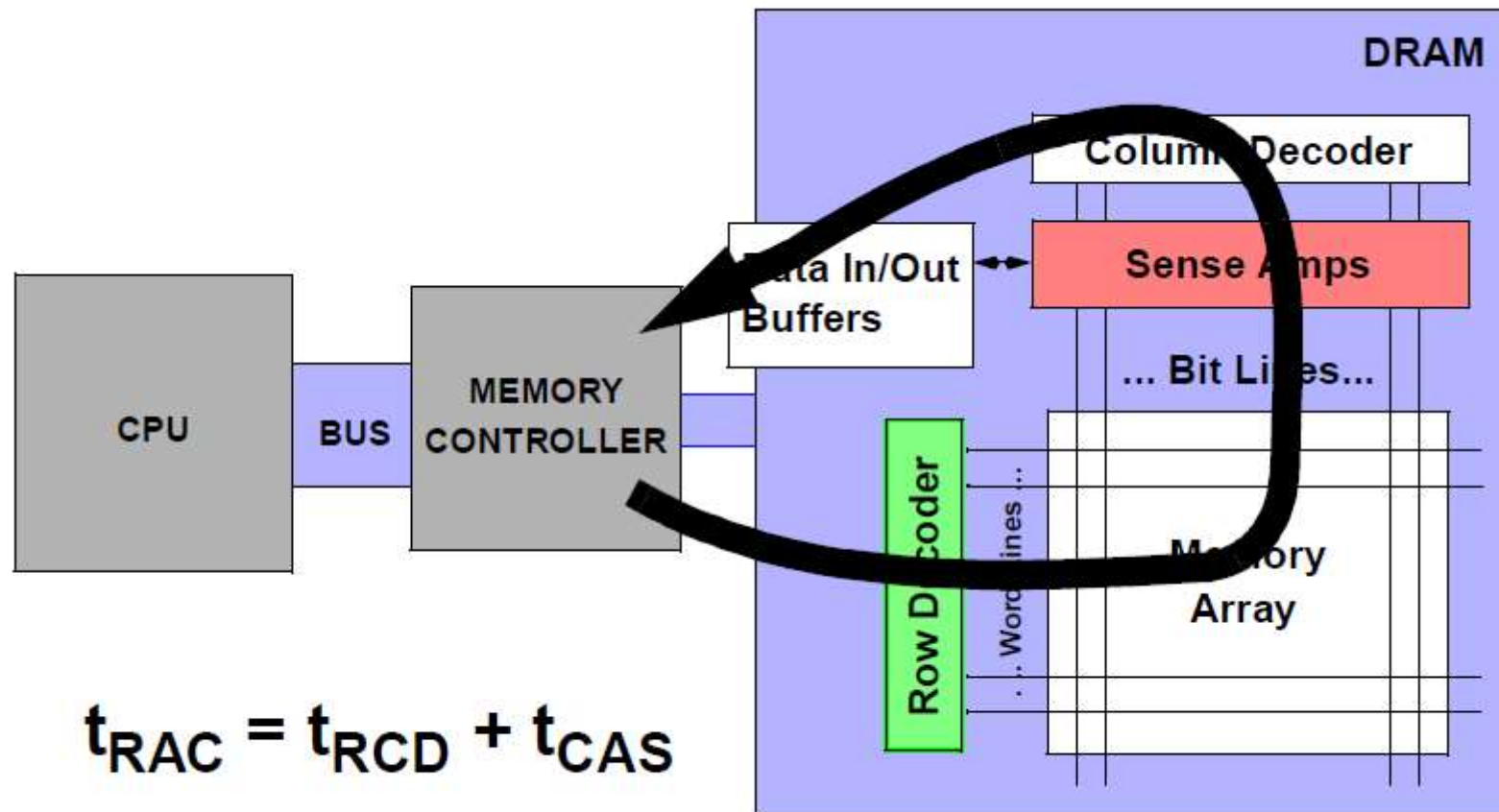
How fast can I get data out of sense amps into memory controller?

27



How fast can I move data from DRAM cell into memory controller?

28



$$t_{RAC} = t_{RCD} + t_{CAS}$$

RAC (Random Access Delay)

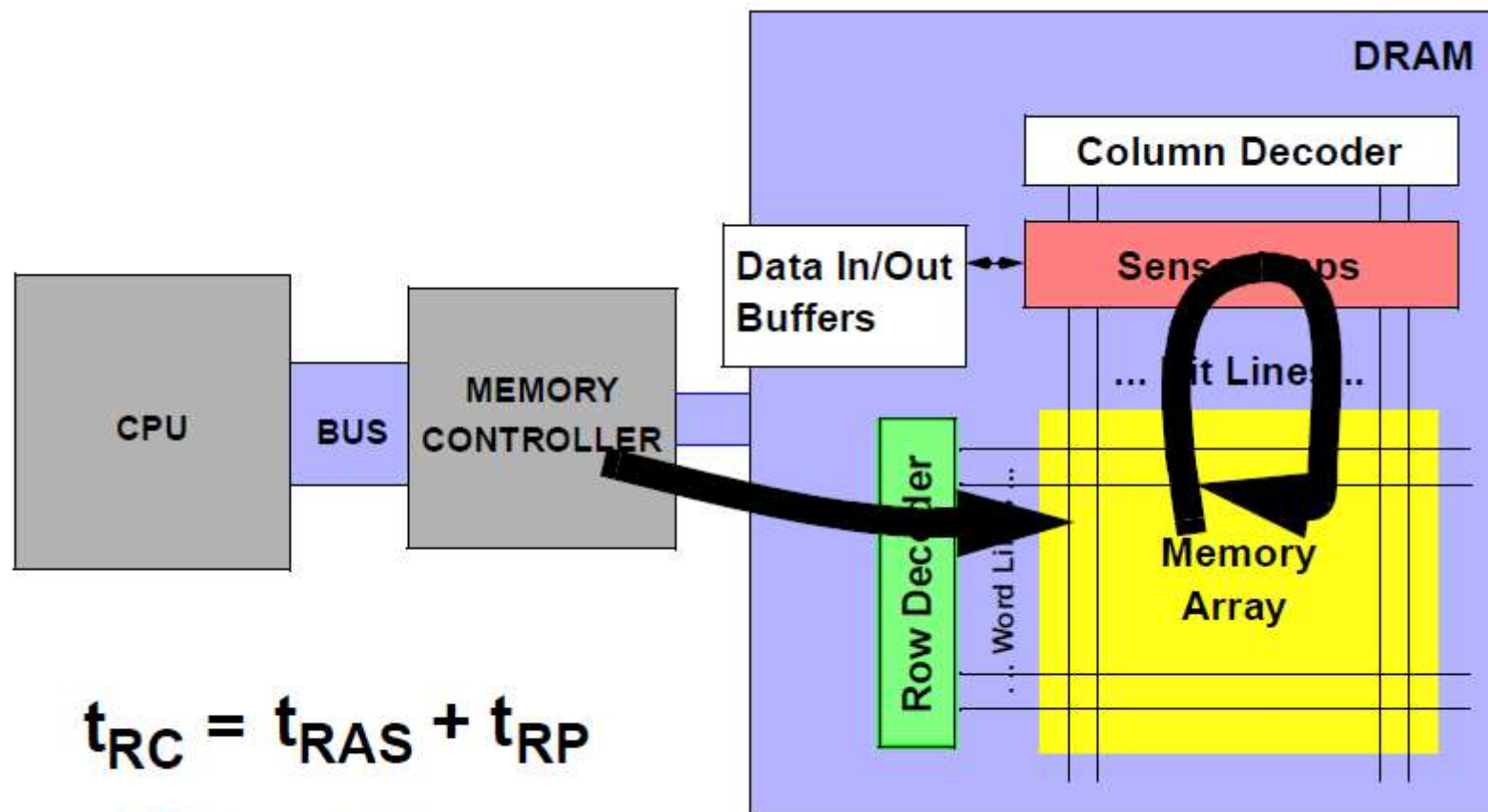
29



VLSI System Design

How fast can I read from different rows?

30

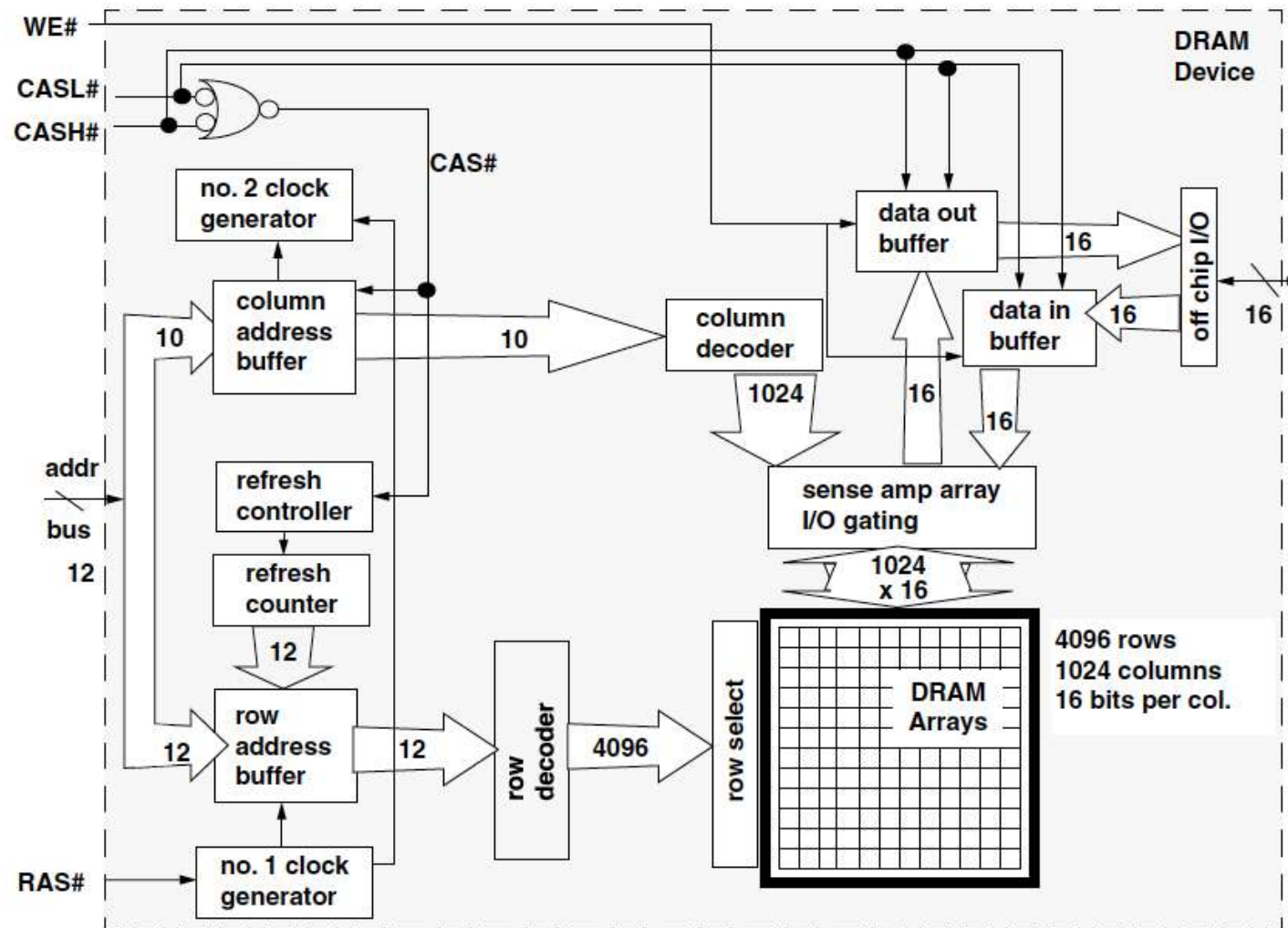


$$t_{RC} = t_{RAS} + t_{RP}$$

RC (Row Cycle Time)

DRAM Organization

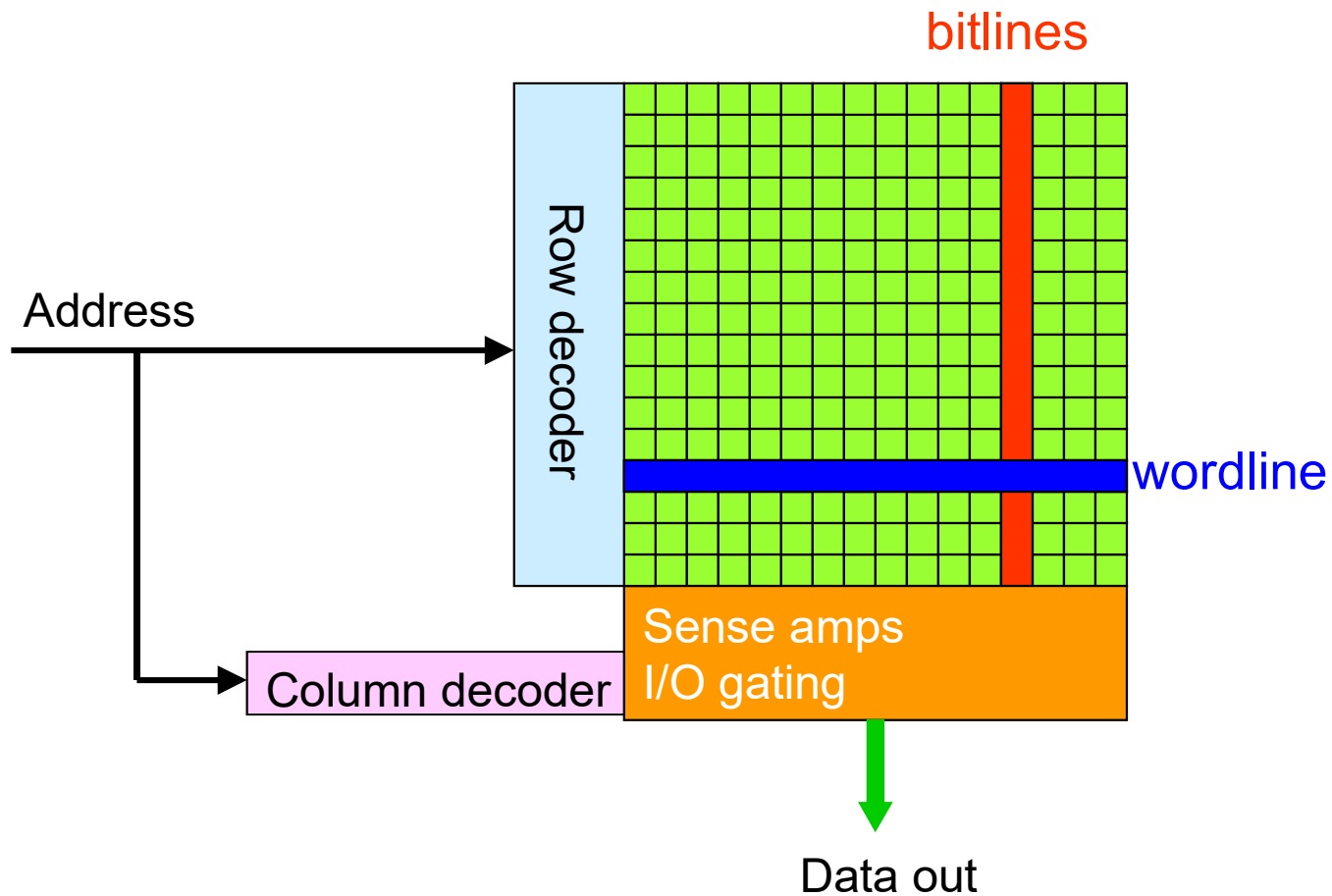
31



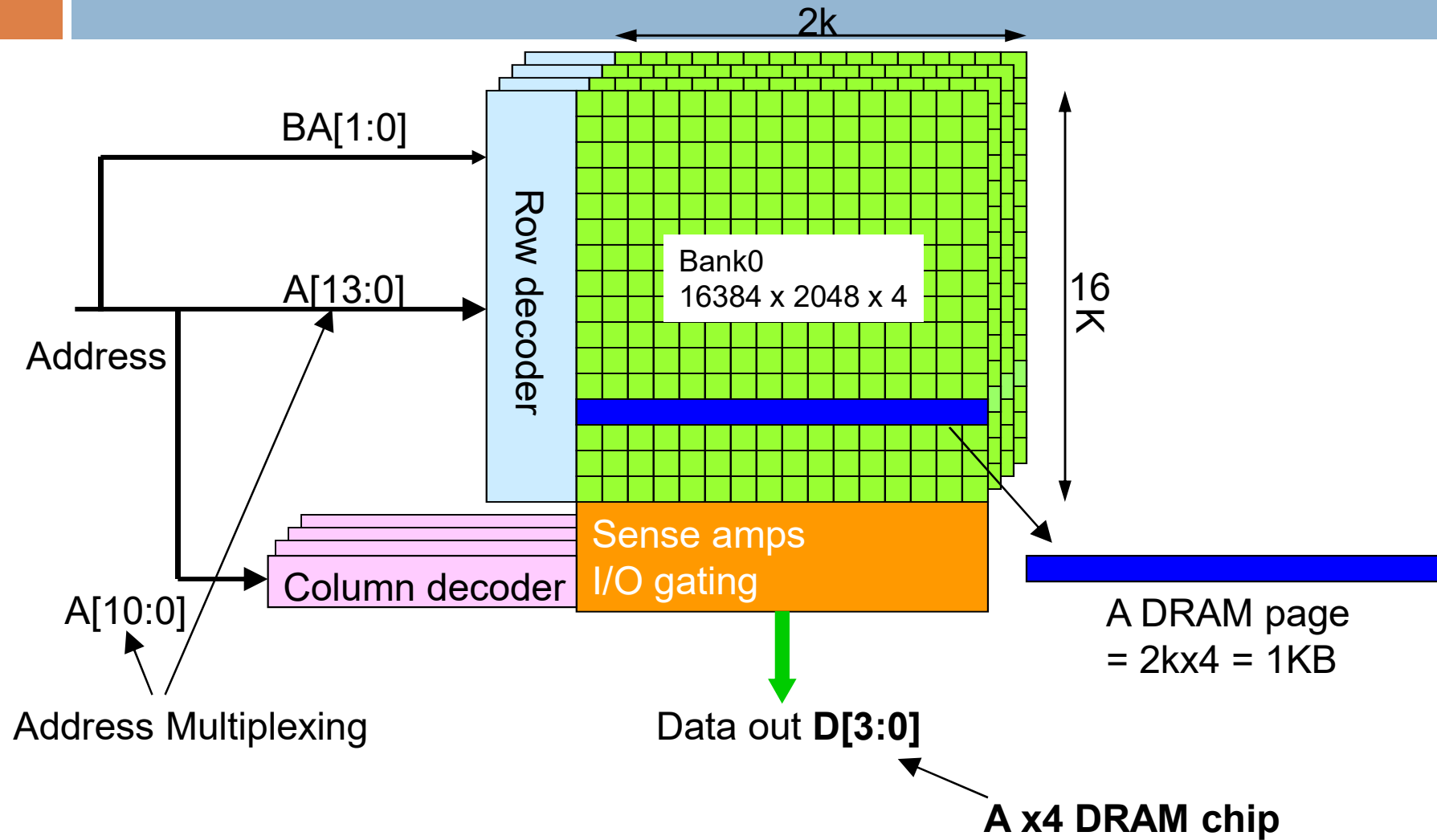
- “Banks” of independent memory arrays inside of a DRAM Chip



One DRAM Bank



Example: 512Mb 4-bank DRAM (x4)



DRAM Refresh

35

- Leaky storage
- Periodic Refresh across DRAM rows
- Un-accessible when refreshing
- Read, and write the same data back

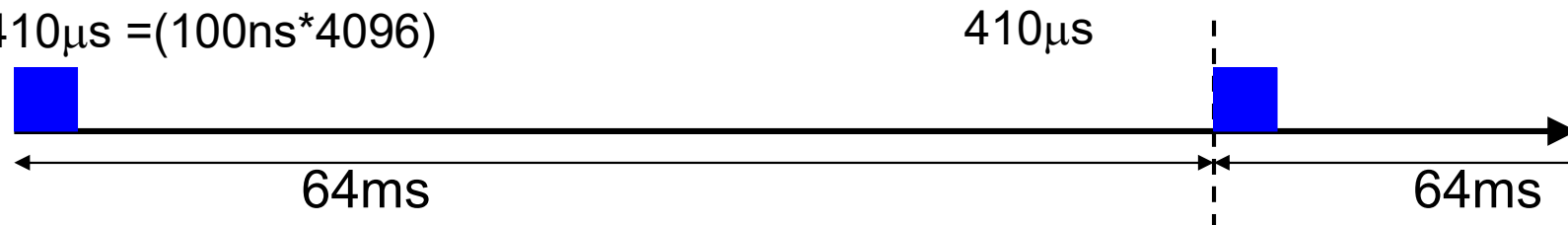
- Example:
 - ▣ 4k rows in a DRAM
 - ▣ 100ns read cycle
 - ▣ Decay in 64ms

 - ▣ $4096 * 100\text{ns} \cong 410\mu\text{s}$ to refresh once
 - ▣ $410\mu\text{s} / 64\text{ms} = 0.64\%$ unavailability

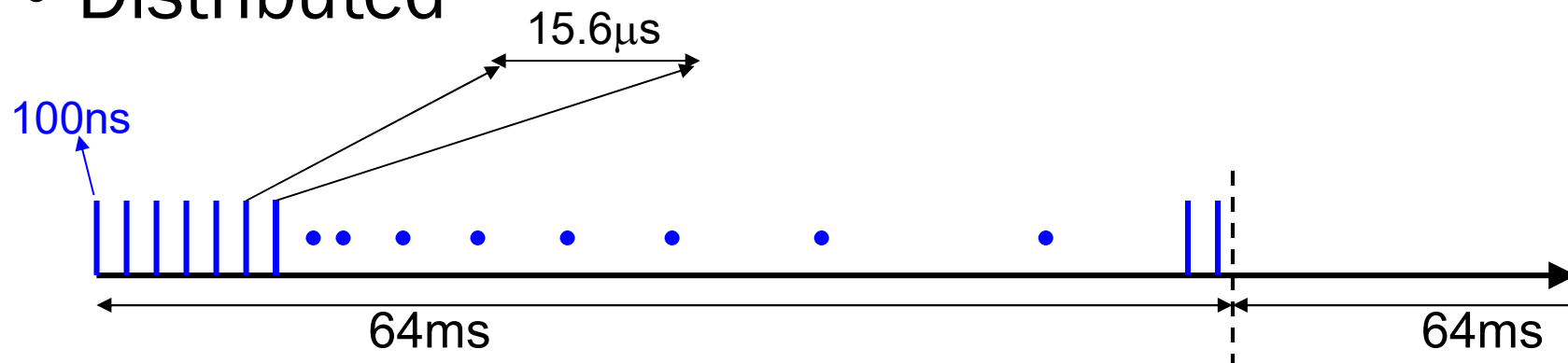
DRAM Refresh Styles

□ Bursty

$$410\mu\text{s} = (100\text{ns} \times 4096)$$



• Distributed



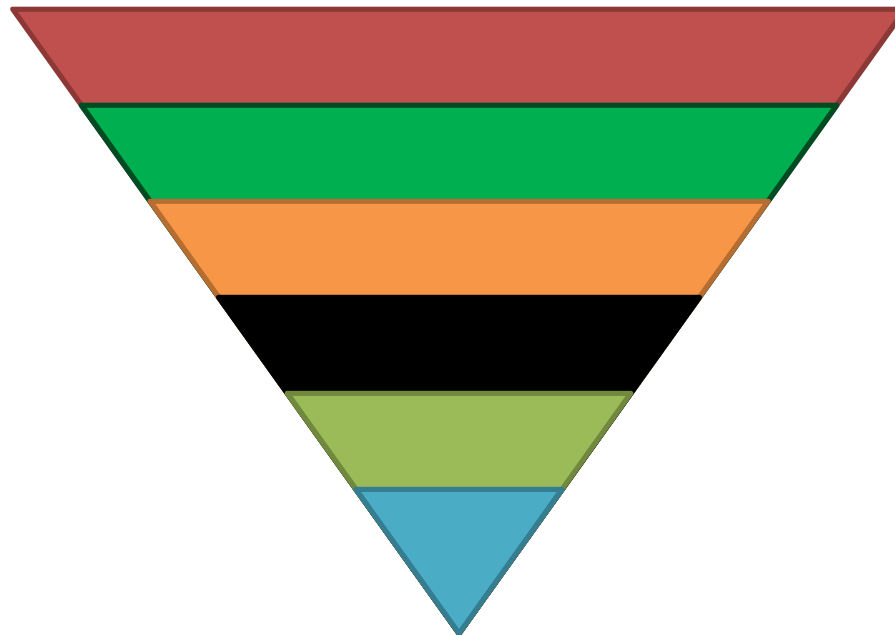
37

DRAM Subsystem Organization

DRAM Subsystem Organization

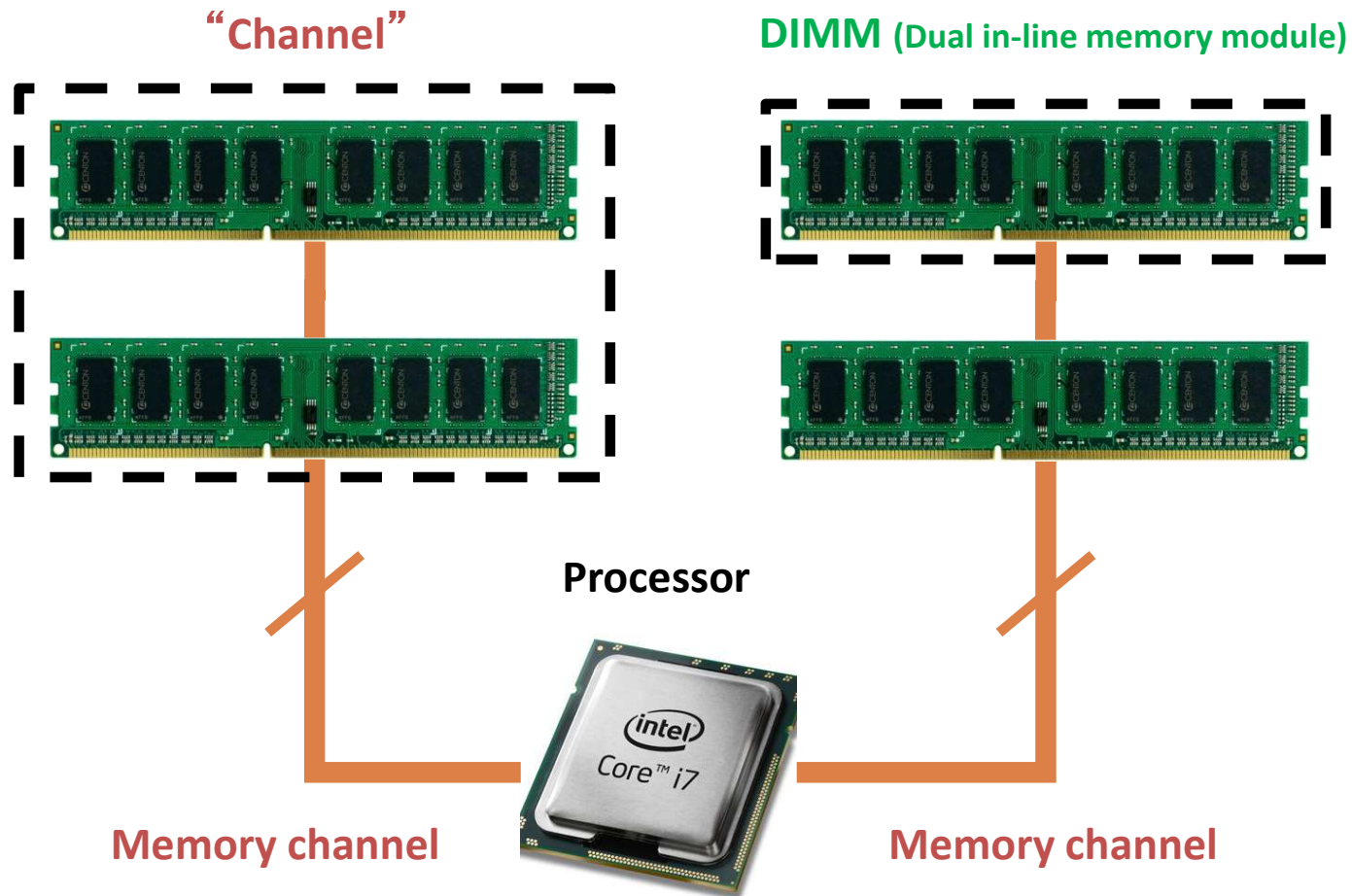
38

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



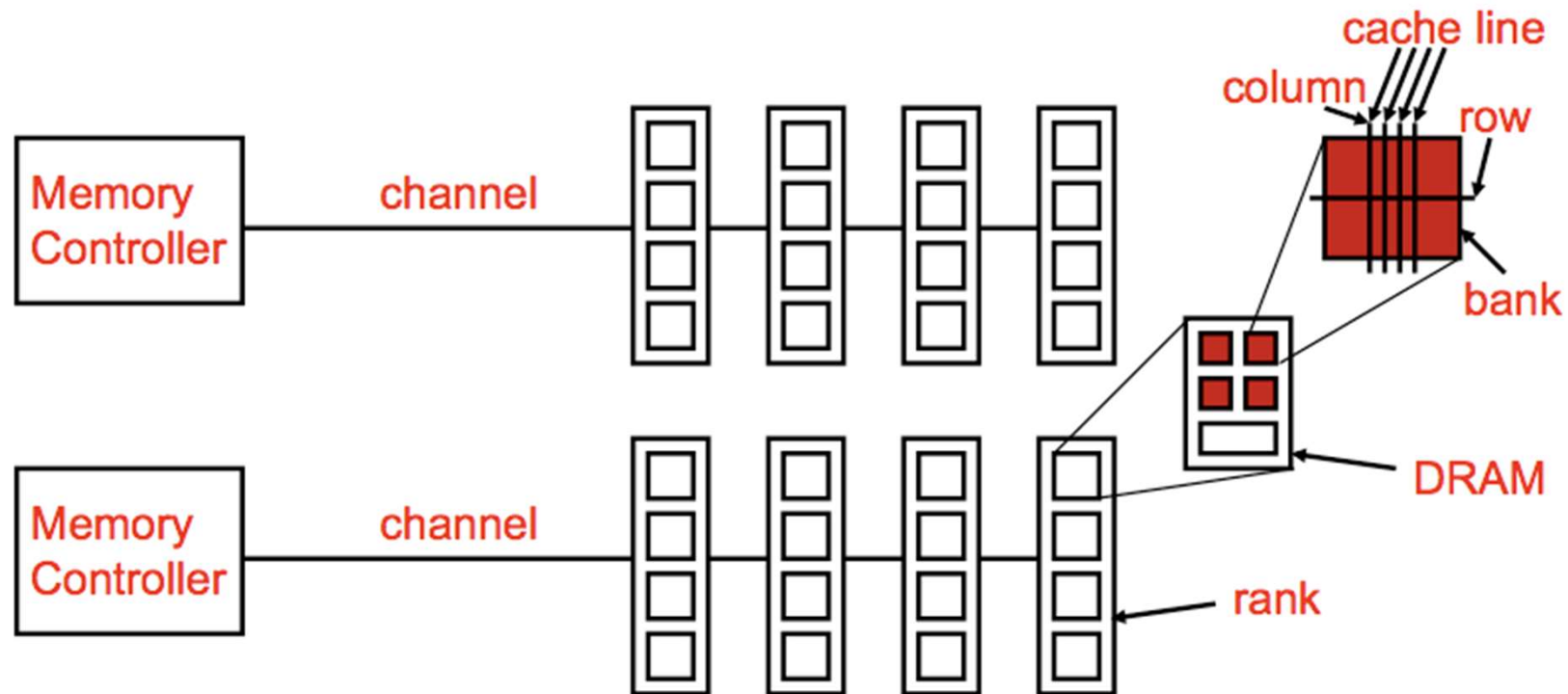
The DRAM subsystem

39



Generalized Memory Structure

40



Rank and Module

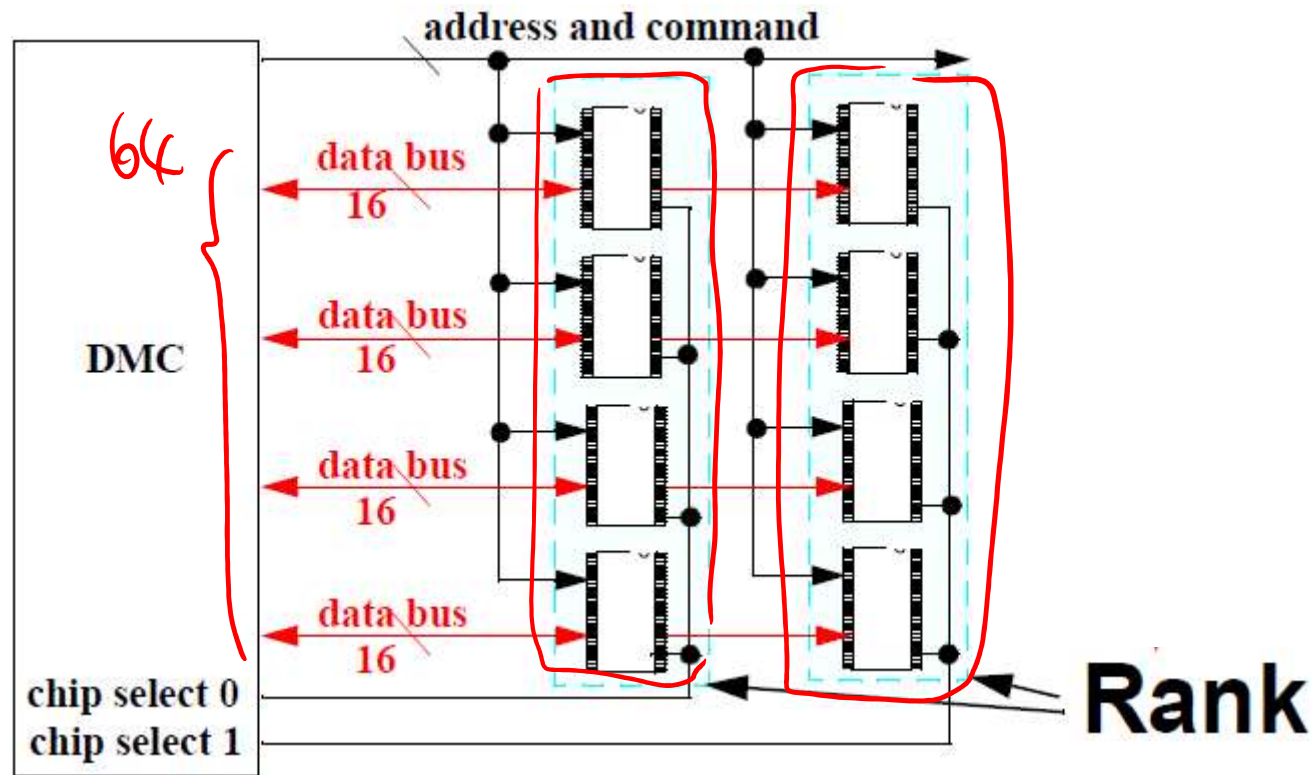
41

- 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

Rank

42

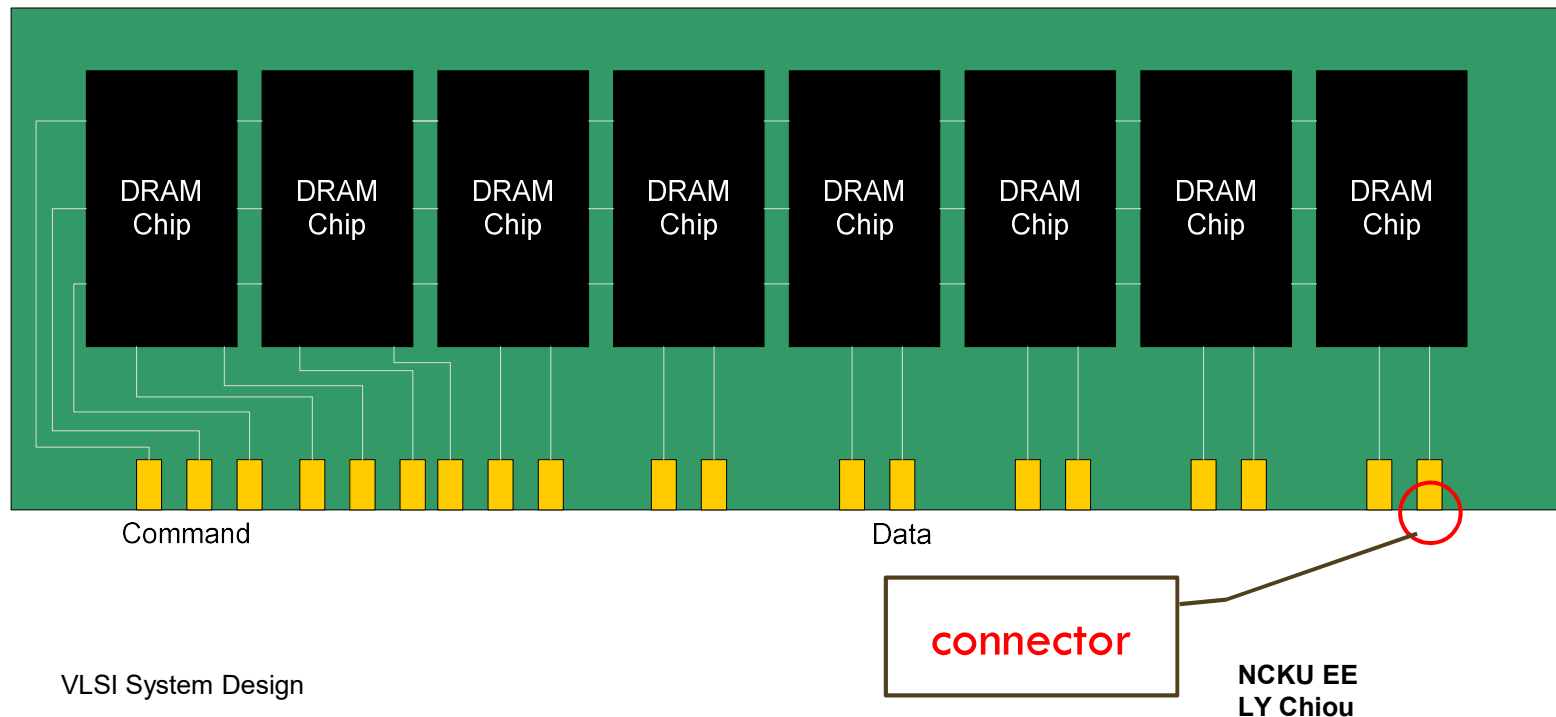
- This example has two ranks, four chips/rank, 16 bit data out/chip



A 64-bit Wide DIMM (One Rank)

43

- Dual-in-line memory module (DIMM), connectors on both are independent, i.e., not connected



Breaking down a DIMM

44

DIMM (Dual in-line memory module)



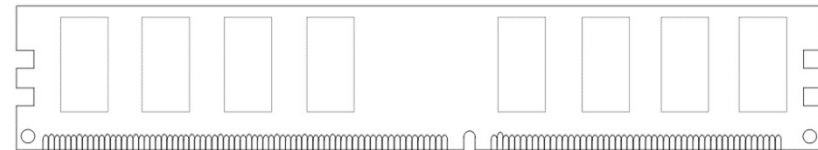
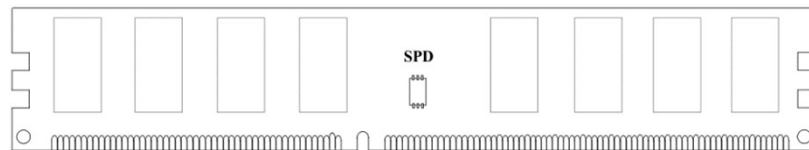
Side view

SIDE

4.00

Front of DIMM

Back of DIMM



Breaking down a DIMM

45

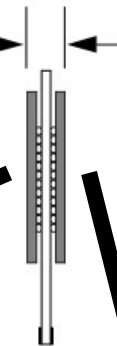
DIMM (Dual in-line memory module)



Side view

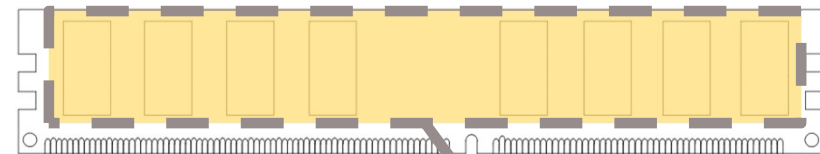
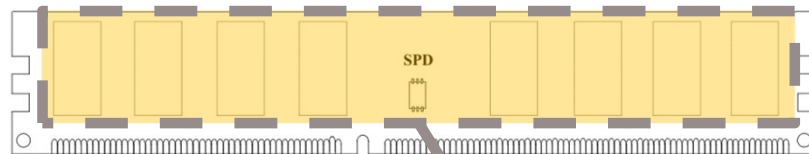
SIDE

4.00



Front of DIMM

Back of DIMM

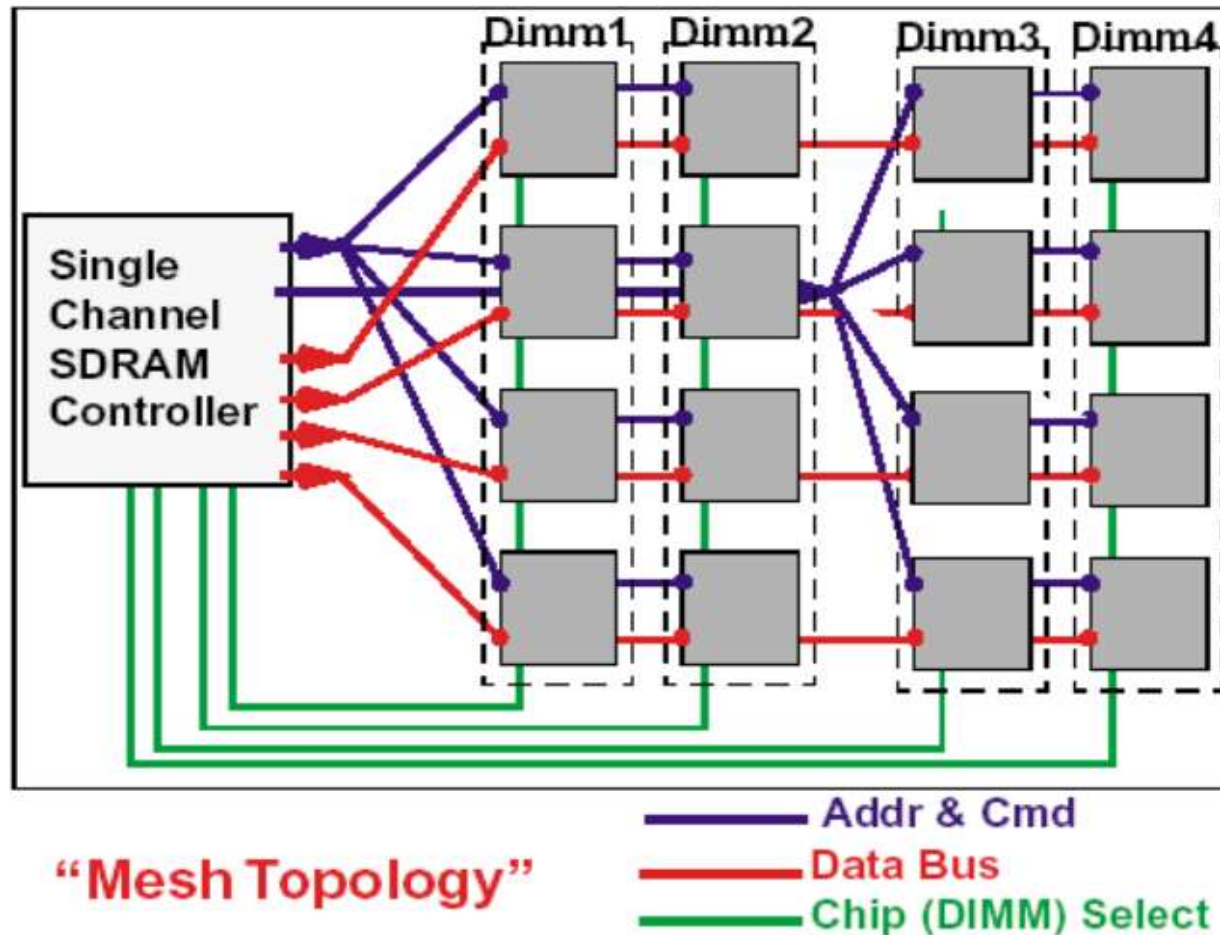


Rank 0: collection of 8 chips

Rank 1

Multiple DIMMs

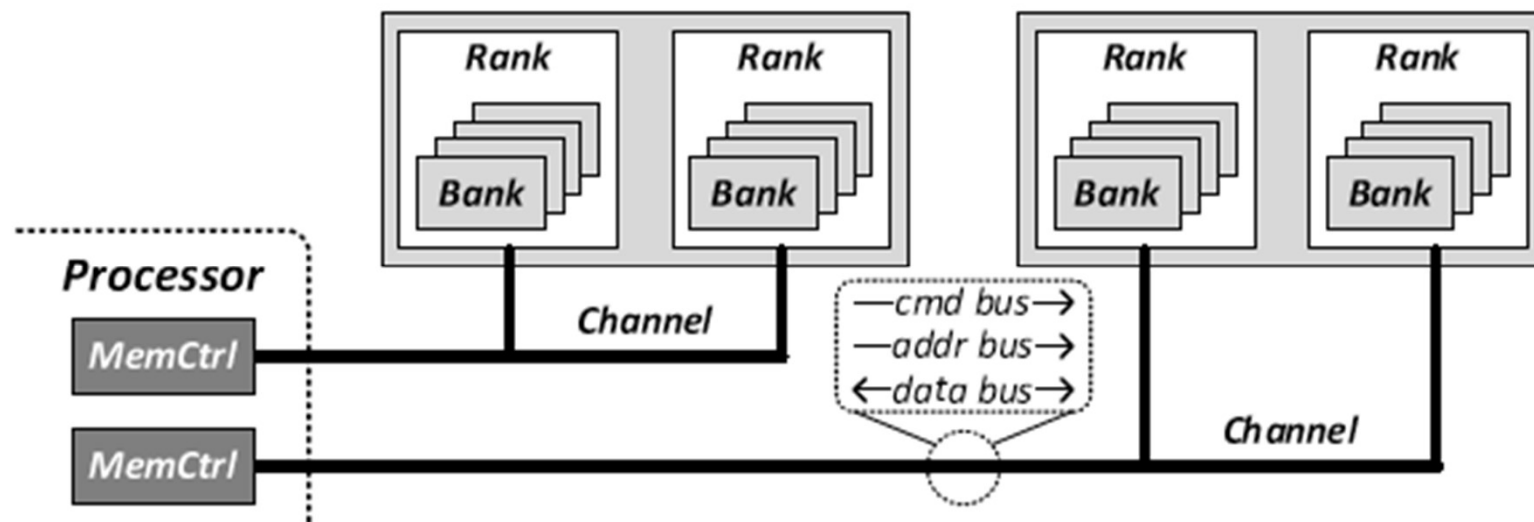
46



- Advantages:
 - ▣ Enables even higher capacity
- Disadvantages:
 - ▣ Interconnect complexity and energy consumption can be high

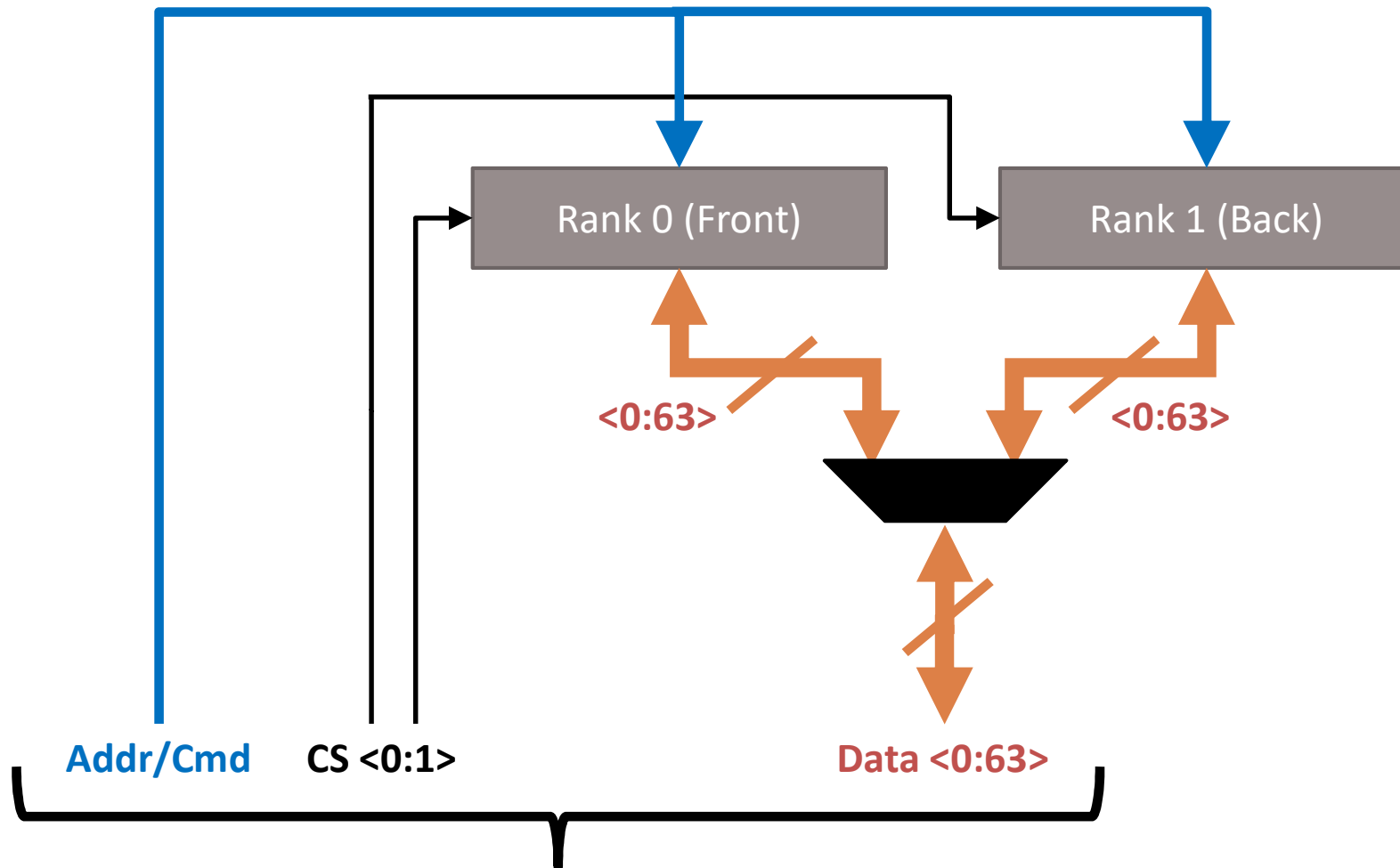
Generalized Memory Structure

47



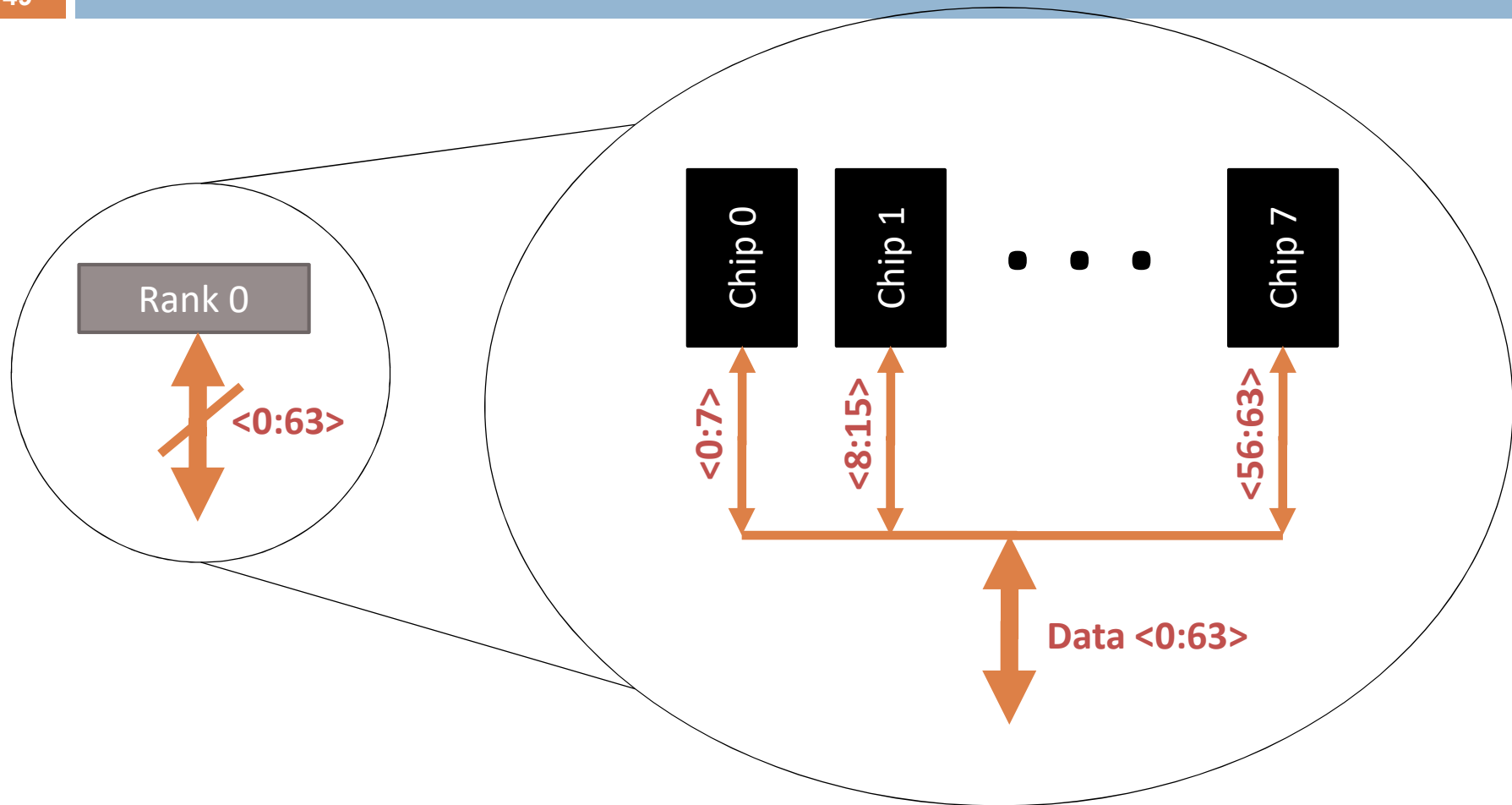
Rank

48



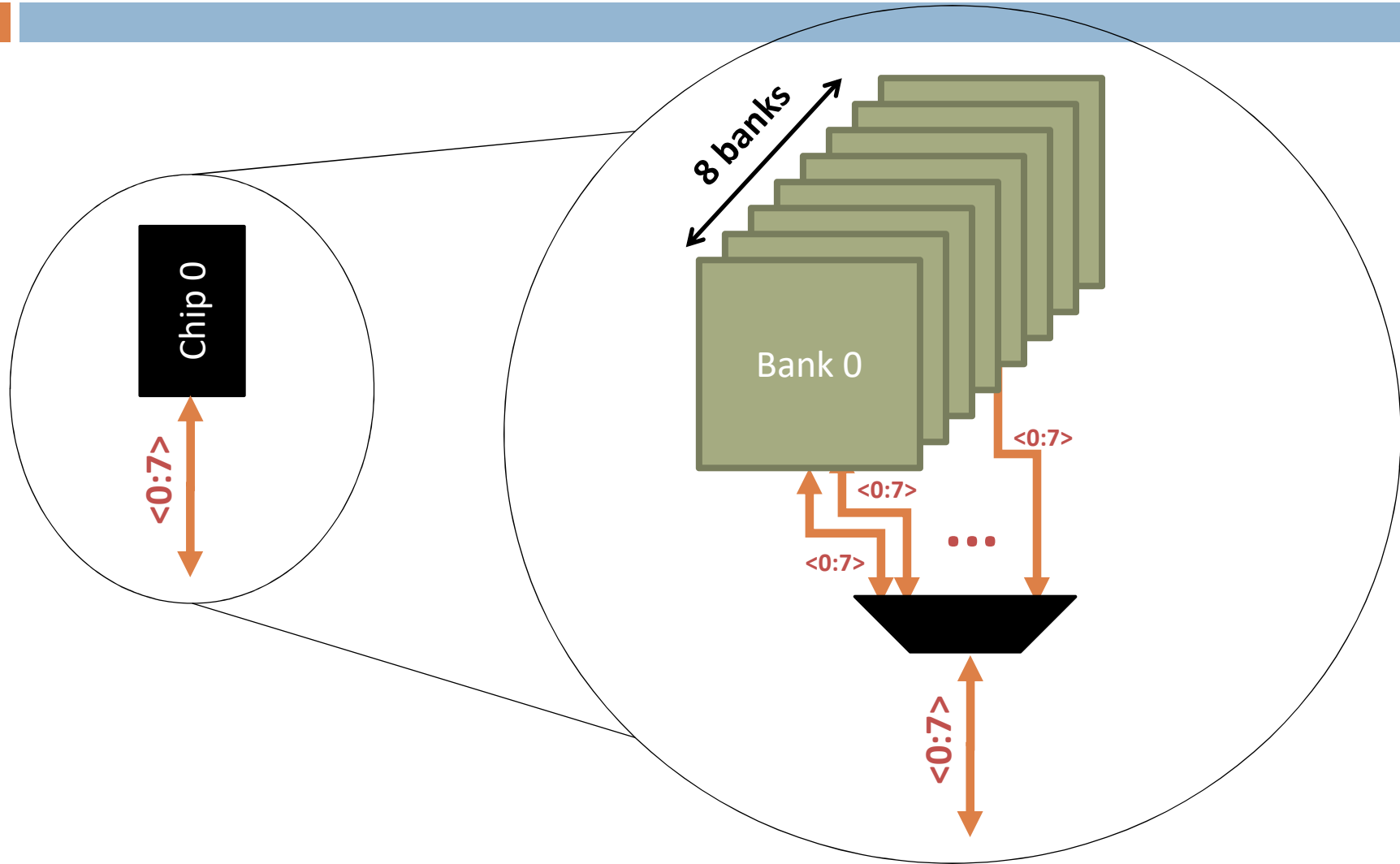
Breaking down a Rank

49



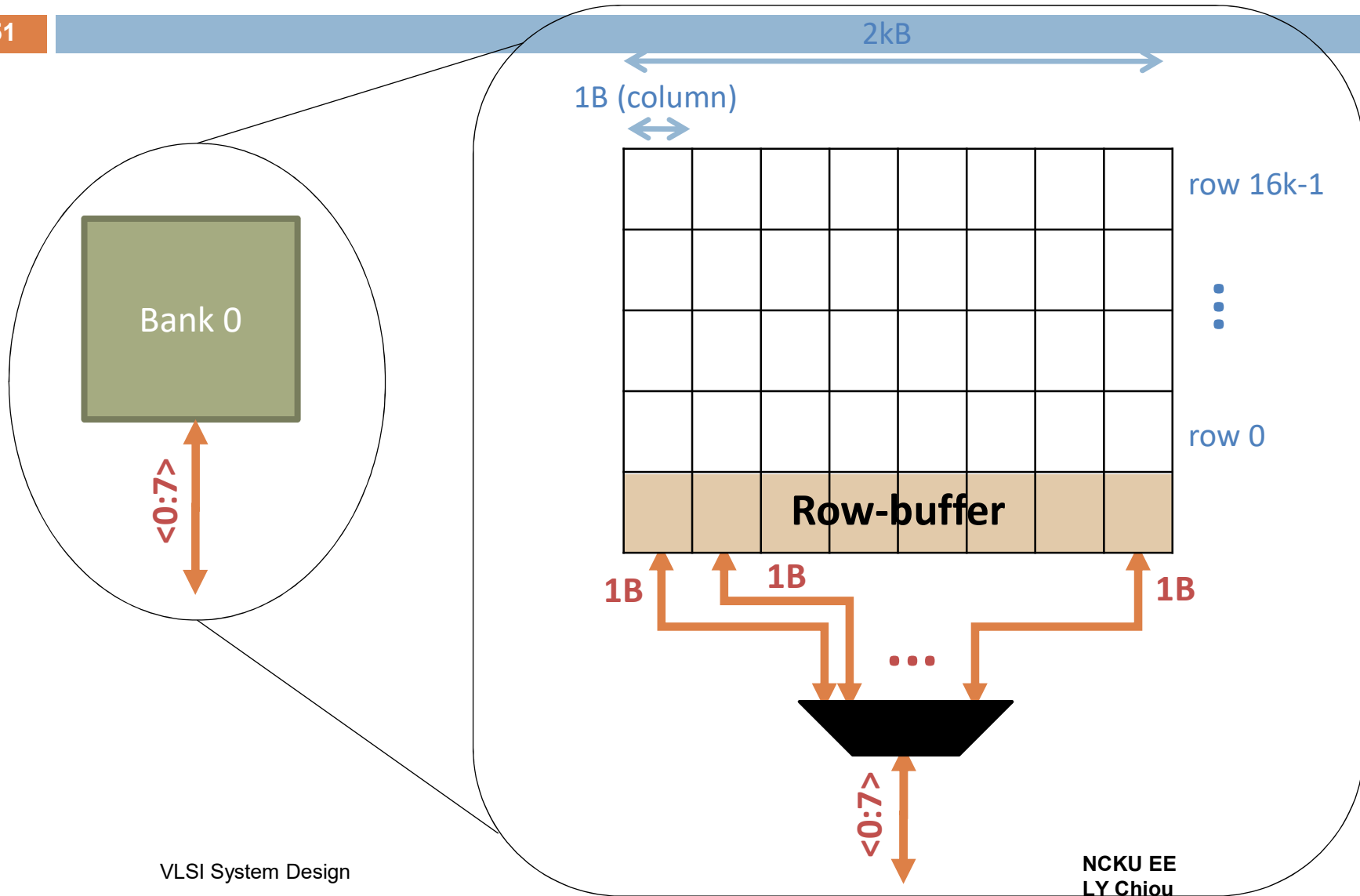
Breaking down a Chip

50



Breaking down a Bank

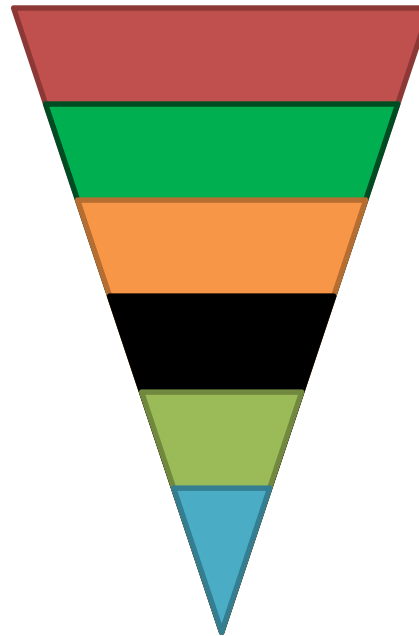
51



DRAM Subsystem Organization

52

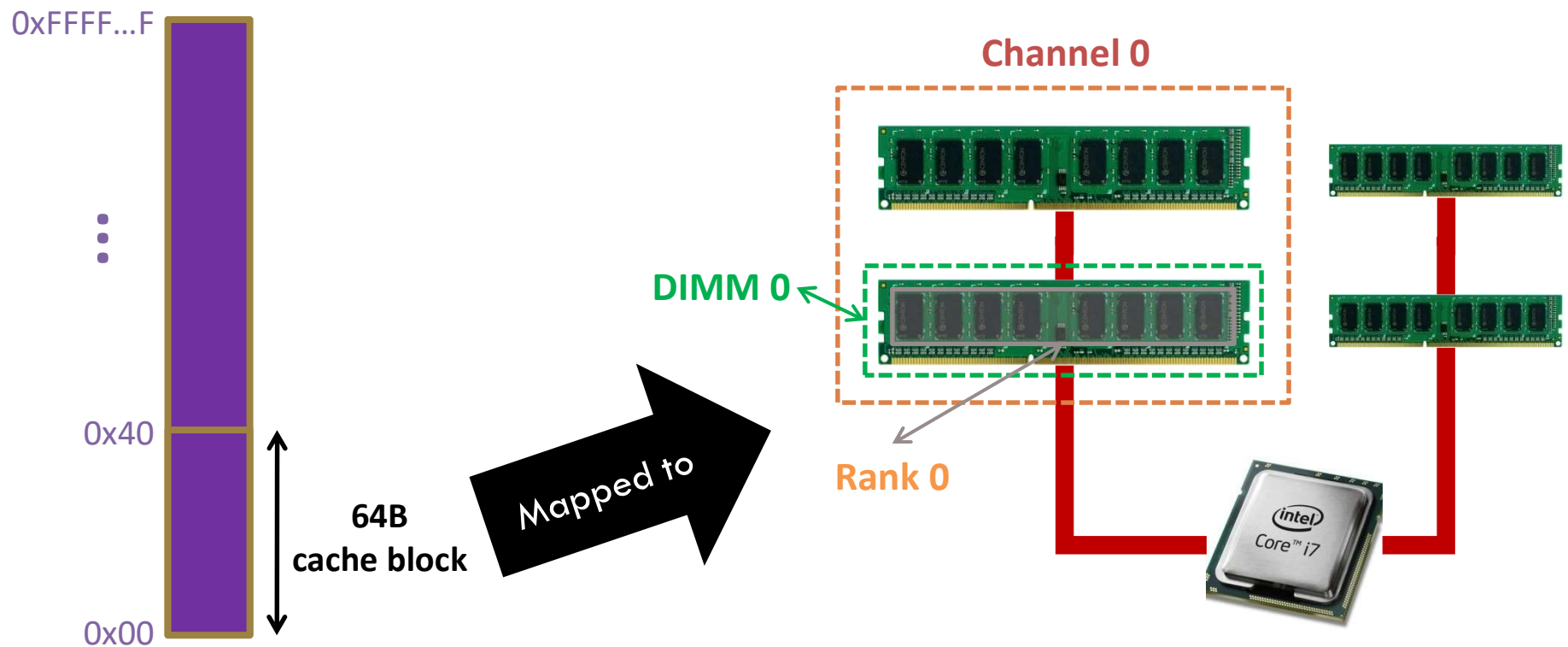
- Channel
- DIMM
- Rank
- Chip
- Bank
- Row/Column



Example: Transferring a cache block

53

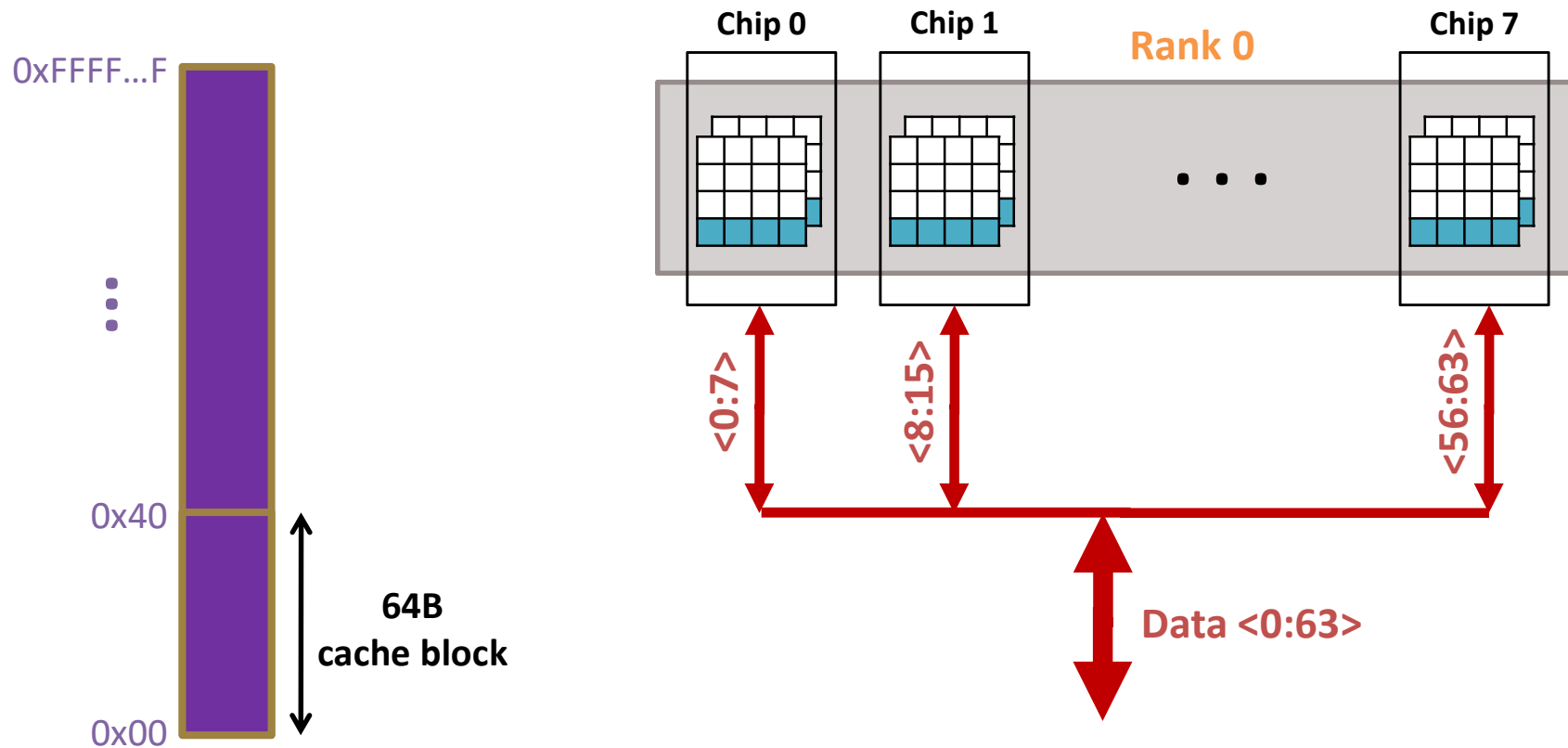
Physical memory space



Example: Transferring a cache block

54

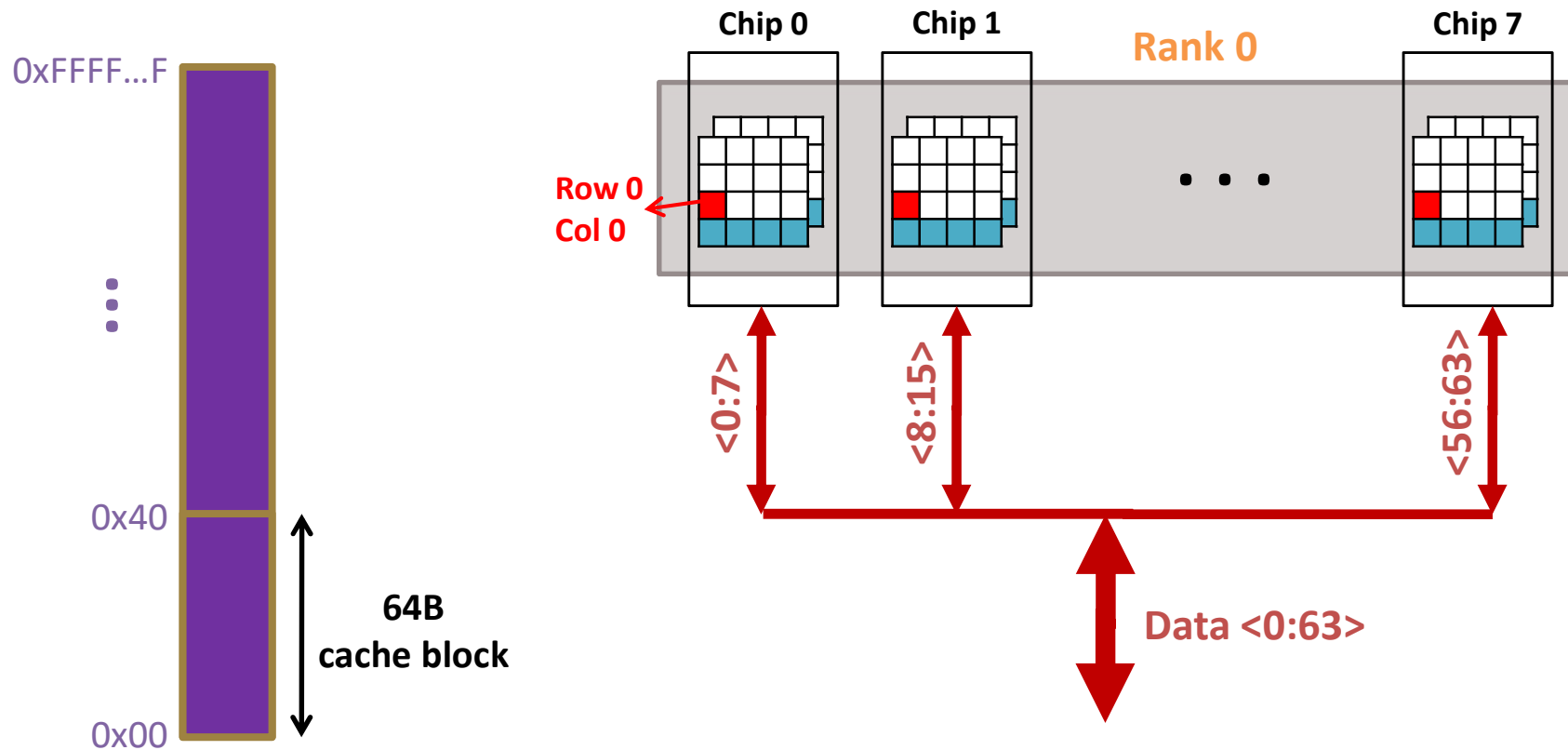
Physical memory space



Example: Transferring a cache block

55

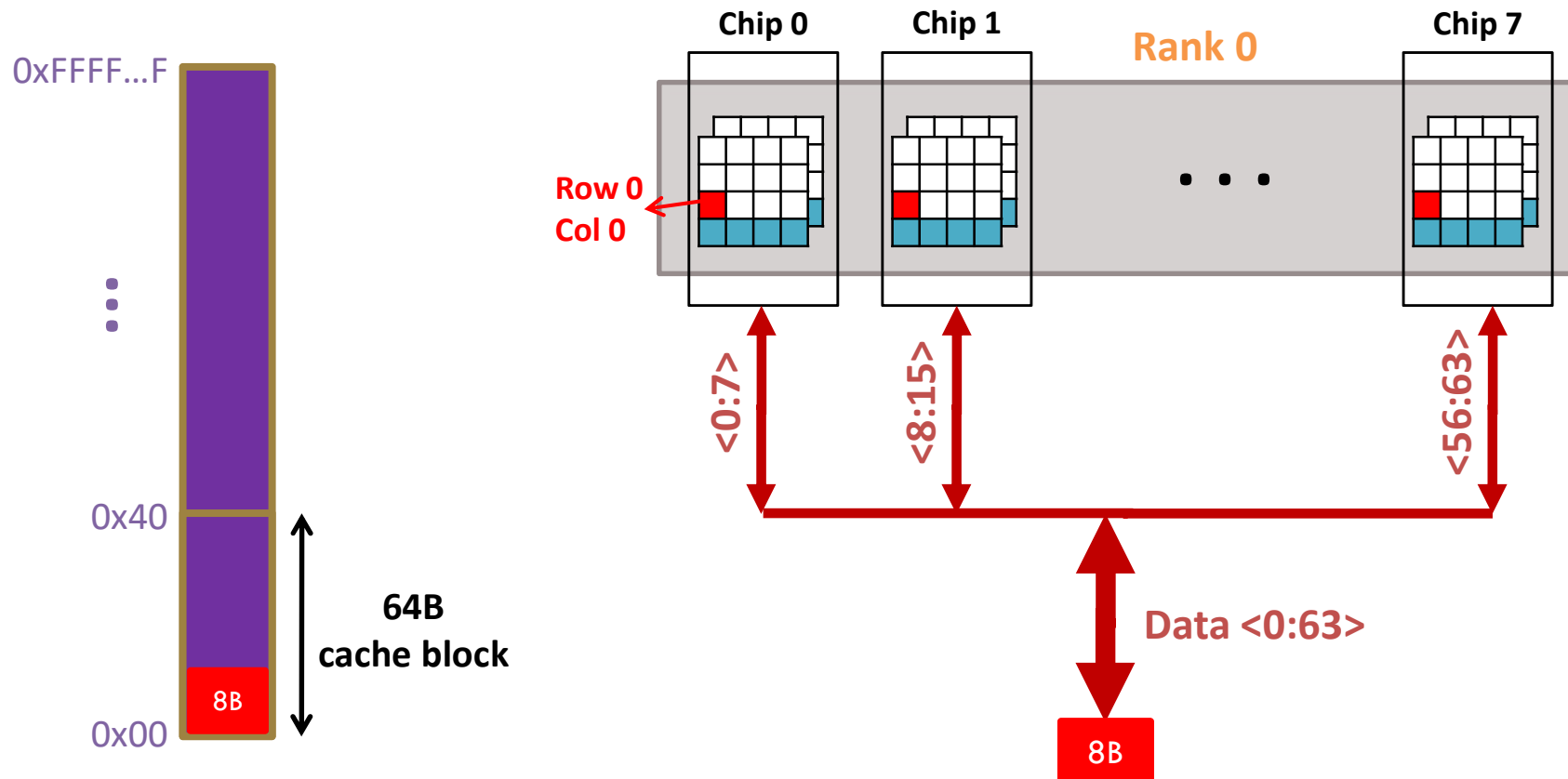
Physical memory space



Example: Transferring a cache block

56

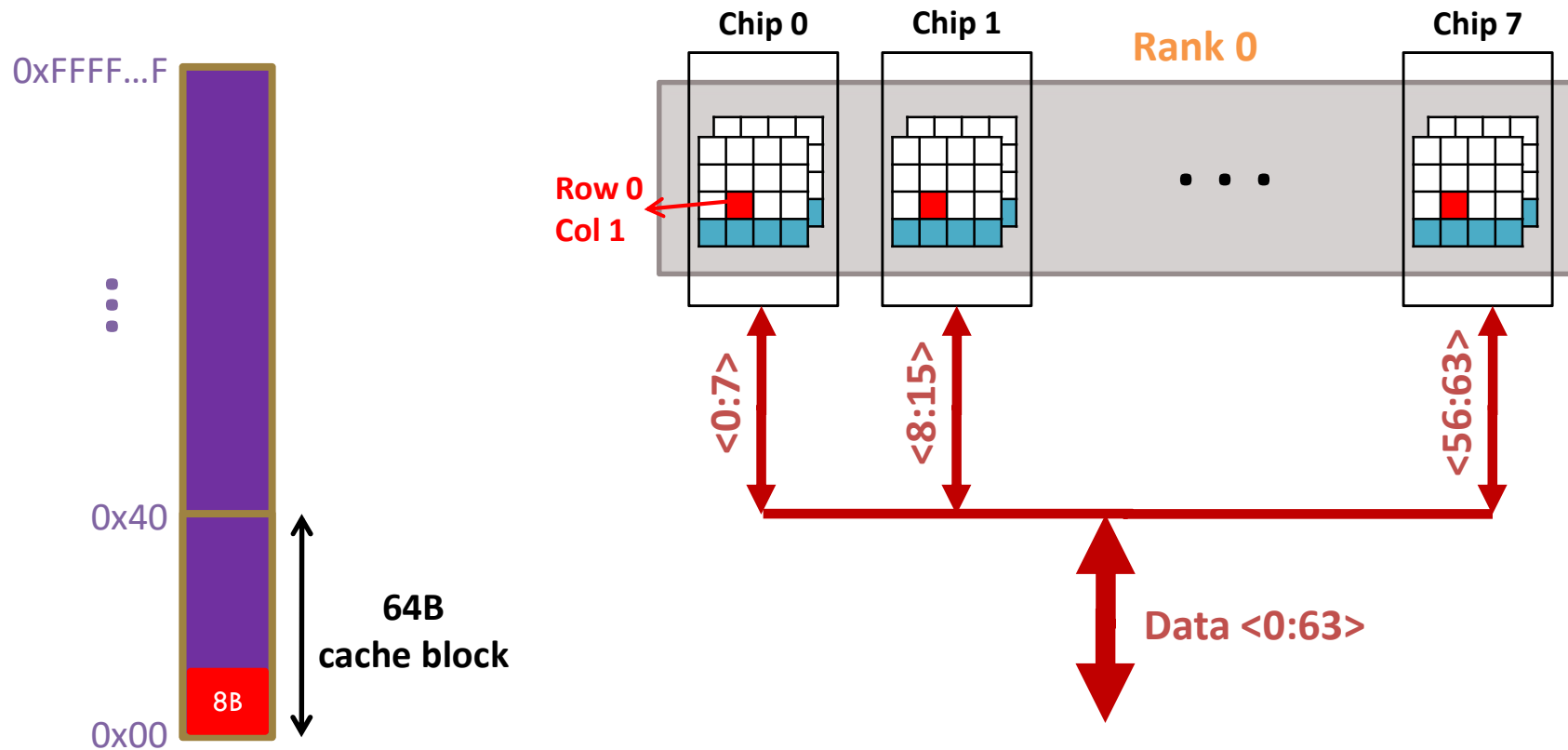
Physical memory space



Example: Transferring a cache block

57

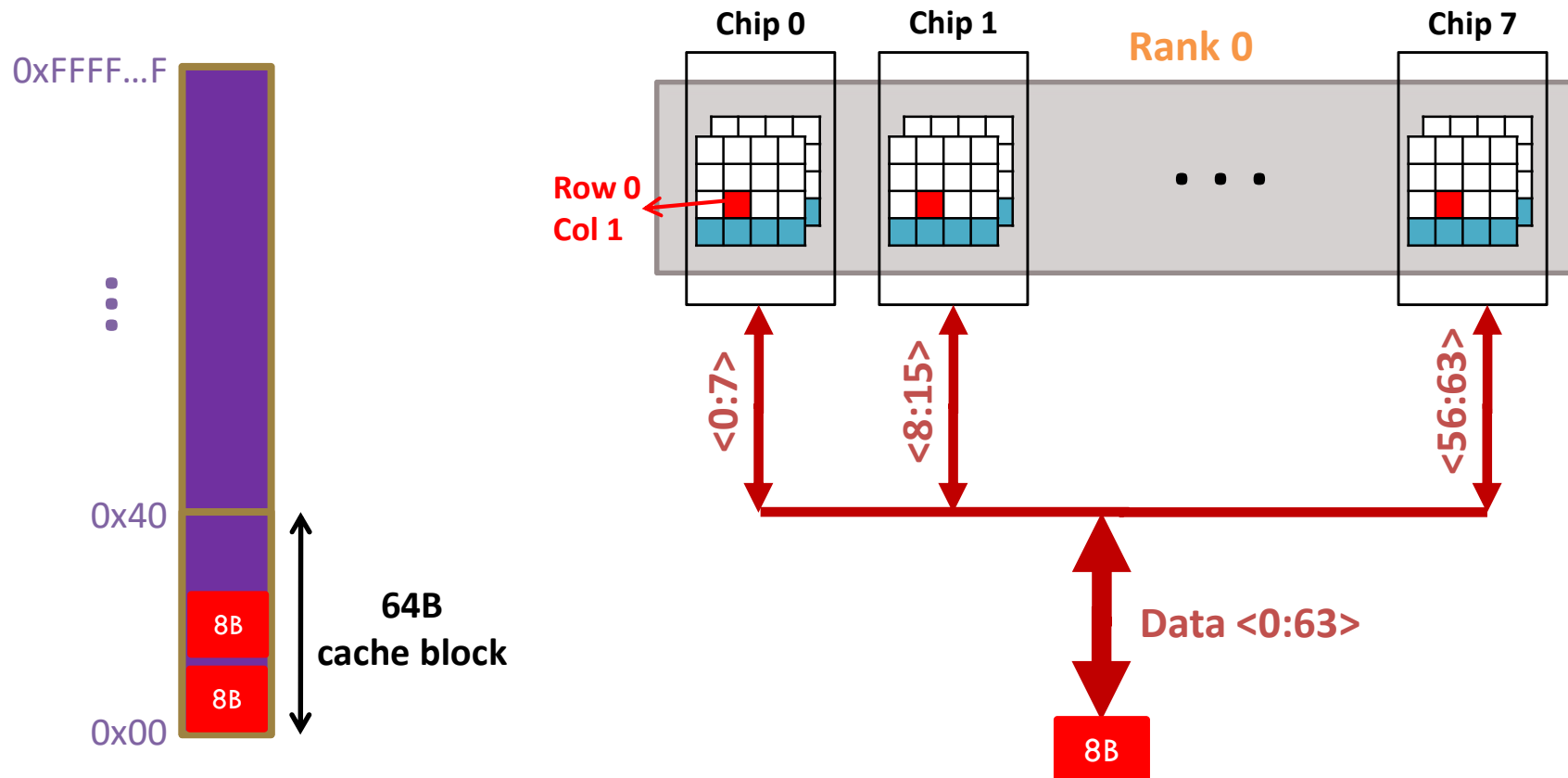
Physical memory space



Example: Transferring a cache block

58

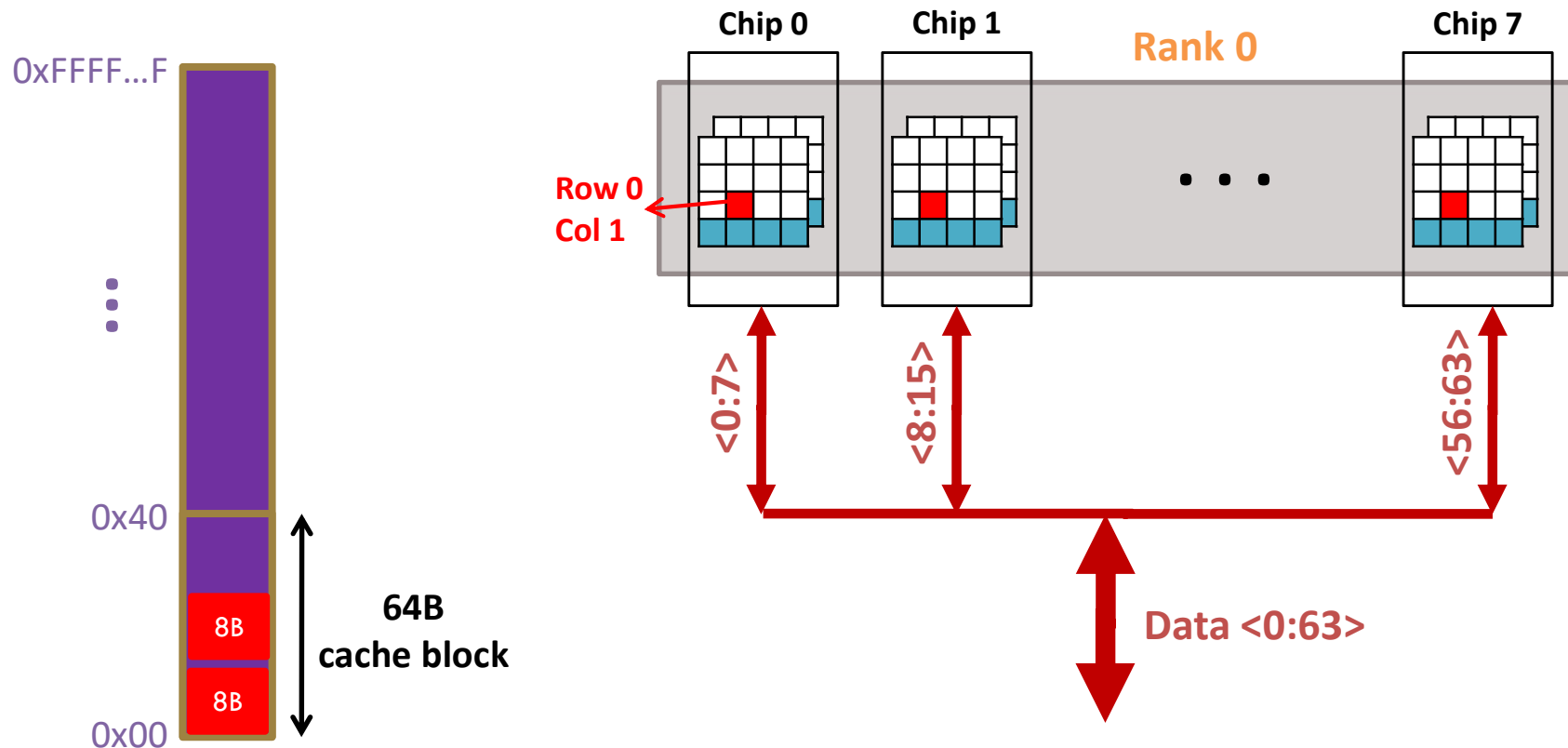
Physical memory space



Example: Transferring a cache block

59

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

60

- 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)

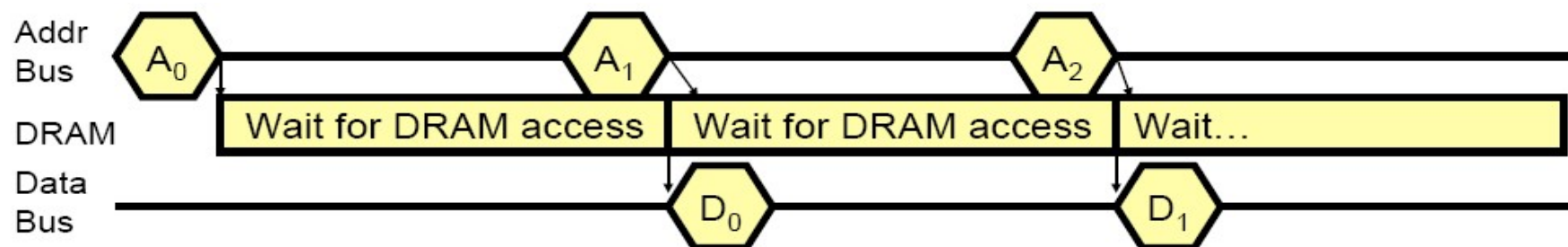
Multiple Banks (Interleaving) and Channels

61

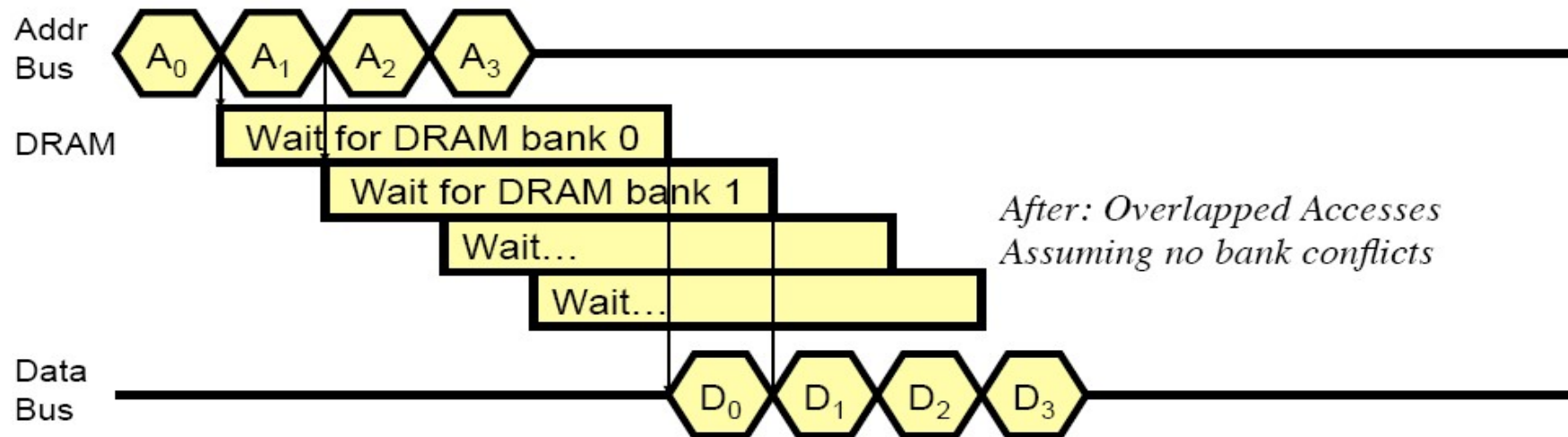
- Multiple banks
 - ▣ Enable **concurrent DRAM accesses**
 - ▣ Bits in address determine which bank an address resides in
- Multiple independent channels serve the same purpose
 - ▣ But they are even better because they have **separate data buses**
 - ▣ **Increased bus bandwidth**
- Enabling more concurrency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

62



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



*After: Overlapped Accesses
Assuming no bank conflicts*

Multiple Channels

63

□ Advantages

- ▣ Increased bandwidth
- ▣ Multiple concurrent accesses (if independent channels)

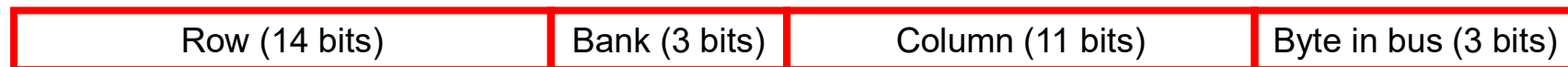
□ Disadvantages

- ▣ Higher cost than a single channel
 - More board wires
 - More pins (if on-chip memory controller)

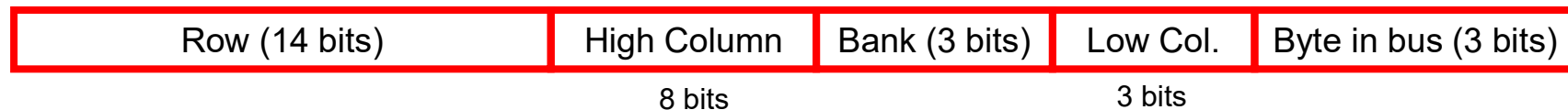
Address Mapping (Single Channel)

64

- 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



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



- Accesses to consecutive cache blocks can be serviced in parallel
- How about random accesses? Strided accesses?

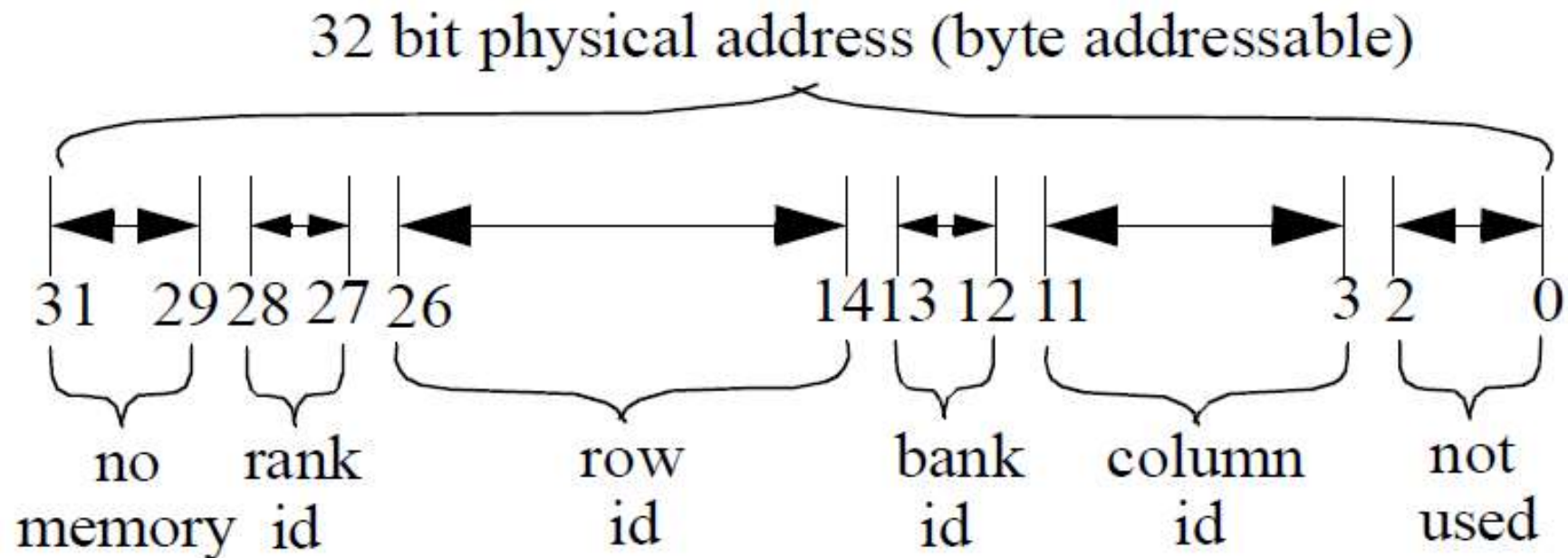
DRAM Basics

65

- Address multiplexing
 - ▣ Send row address when RAS asserted
 - ▣ Send column address when CAS asserted
- DRAM reads are self-destructive
 - ▣ Rewrite after a read
- Memory array
 - ▣ All bits within an array work in unison
- Memory bank
 - ▣ Different banks can operate independently
- DRAM rank
 - ▣ Chips inside the same rank are accessed simultaneously

DRAM Address Structure

66



Notes

67

- The memory controller schedules memory accesses to maximize row buffer hit rates and bank/rank parallelism
- Banks and ranks offer memory parallelism
- Row buffers act as a cache within DRAM
 - ▣ Row buffer hit: ~ 20 ns access time (must only move data from row buffer to pins)
 - ▣ Empty row buffer access: ~ 40 ns (must first read arrays, then move data from row buffer to pins)
 - ▣ Row buffer conflict: ~ 60 ns (must first writeback the existing row, then read new row, then move data to pins)
- In addition, must wait in the queue (tens of nano-seconds) and incur address/cmd/data transfer delays (~ 10 ns)

Latency and Power Wall

68

- Both improved by employing smaller arrays
 - ▣ Penalty in density and cost
- Both improved by increasing the row buffer hit rate
 - ▣ Requires intelligent mapping of data to rows, clever scheduling of requests, etc.

69

DRAM System Signaling and Timing

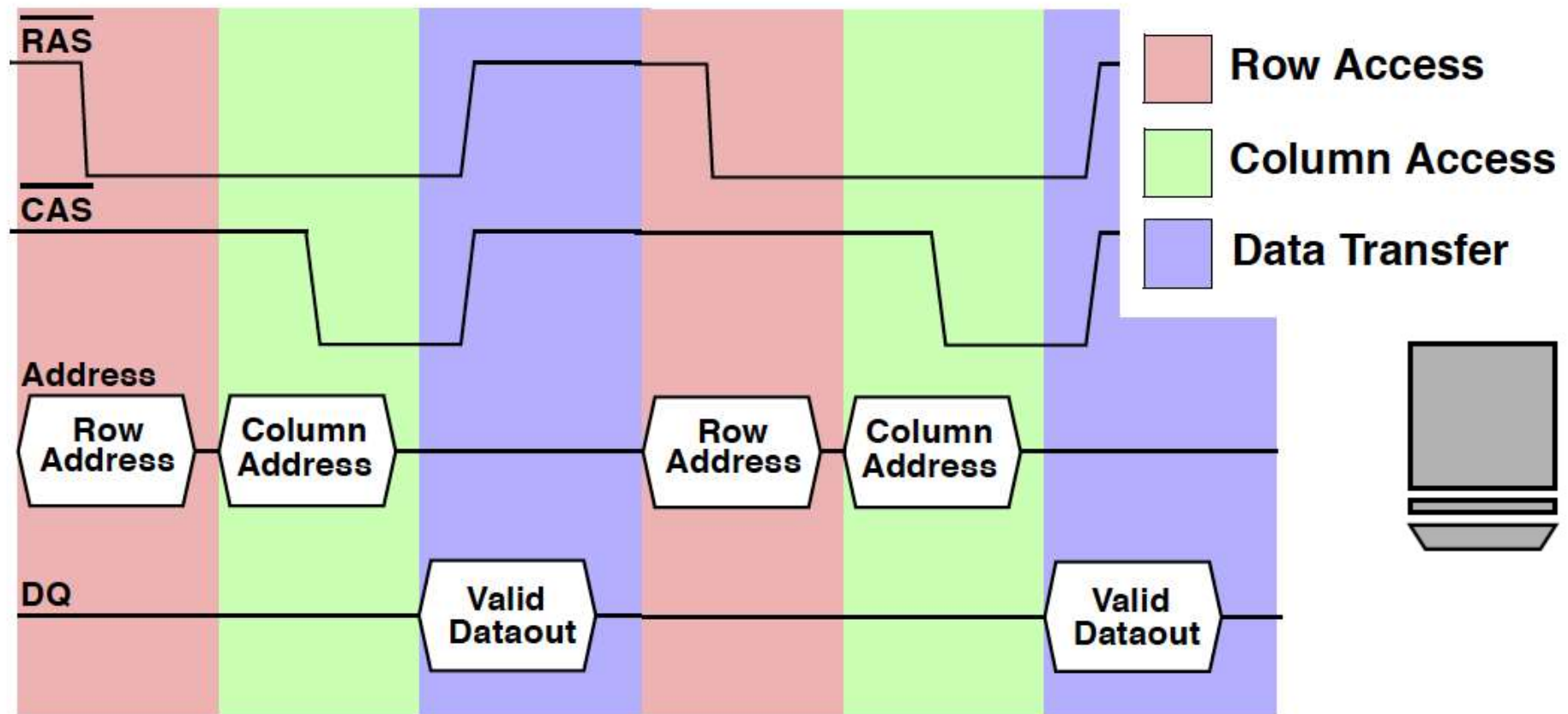
Different DRAM Systems

70

- Innovation targeted towards higher bandwidth for memory systems:
 - SDRAM - synchronous DRAM
 - RDRAM - Rambus DRAM
 - EDORAM - extended data out SRAM
 - Three-dimensional RAM
 - Hyper-page mode DRAM video RAM
 - Multibank DRAM

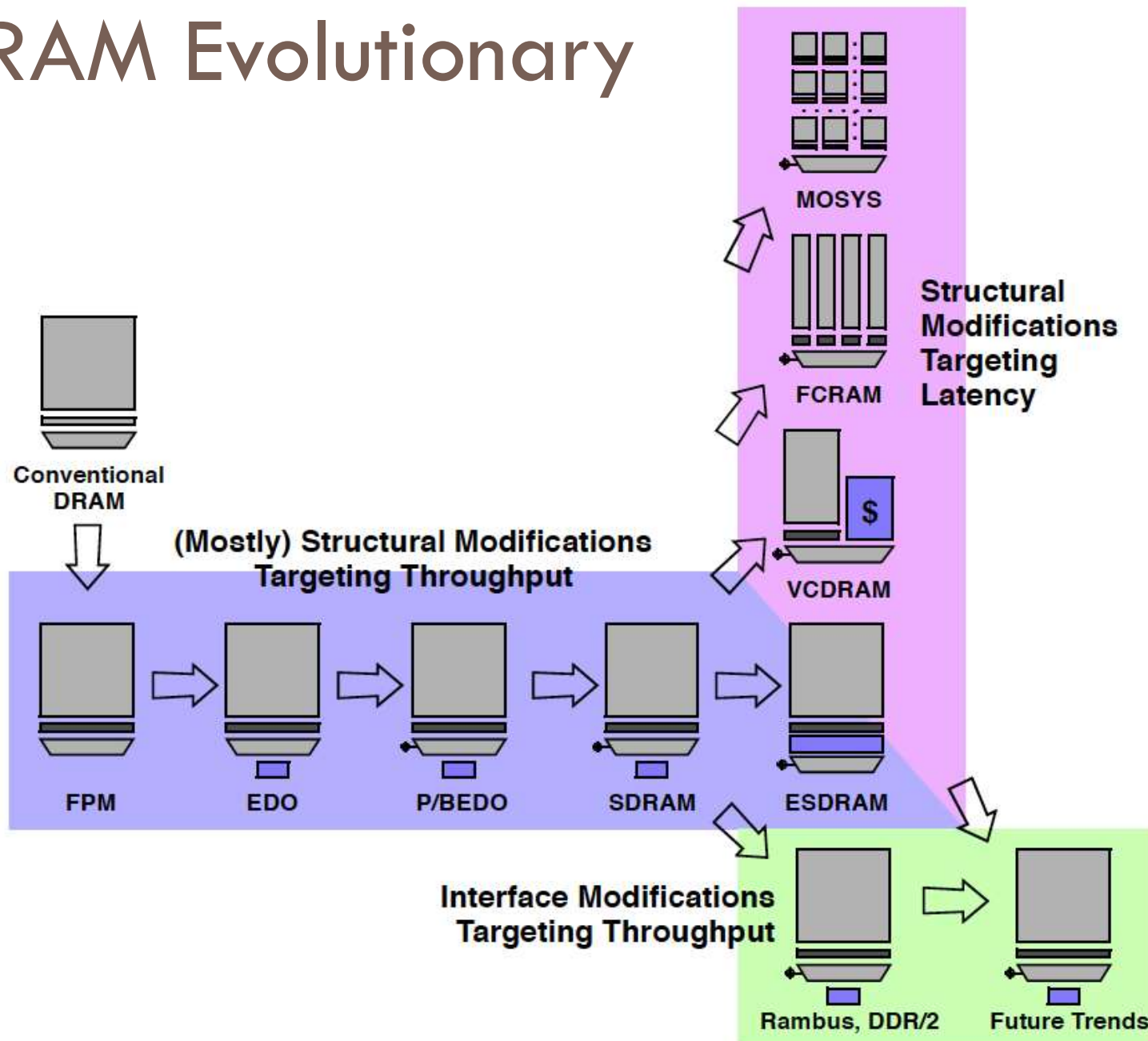
Read Timing of Conventional DRAM

71



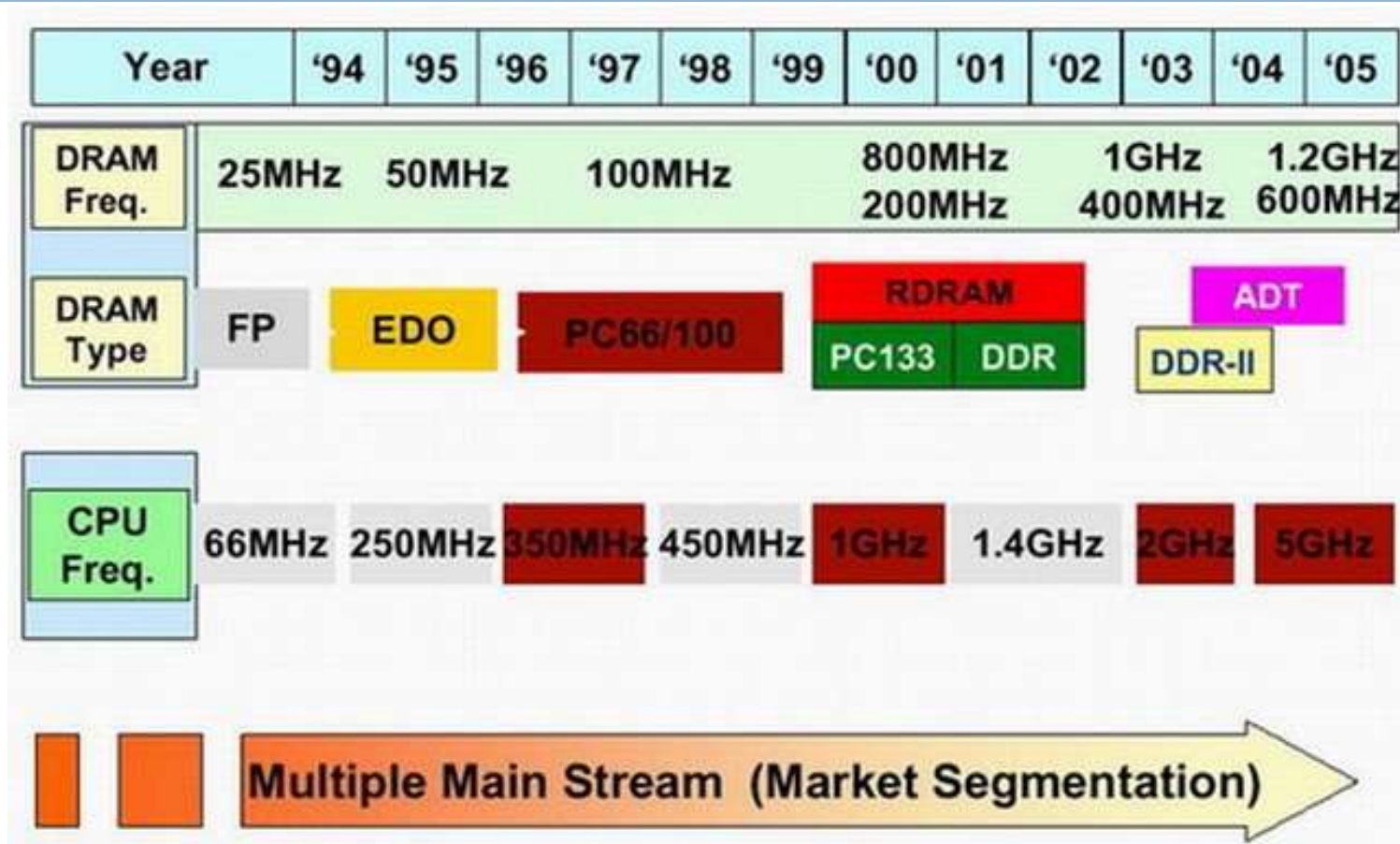
DRAM Evolutionary

72



Frequency of DRAM Generations

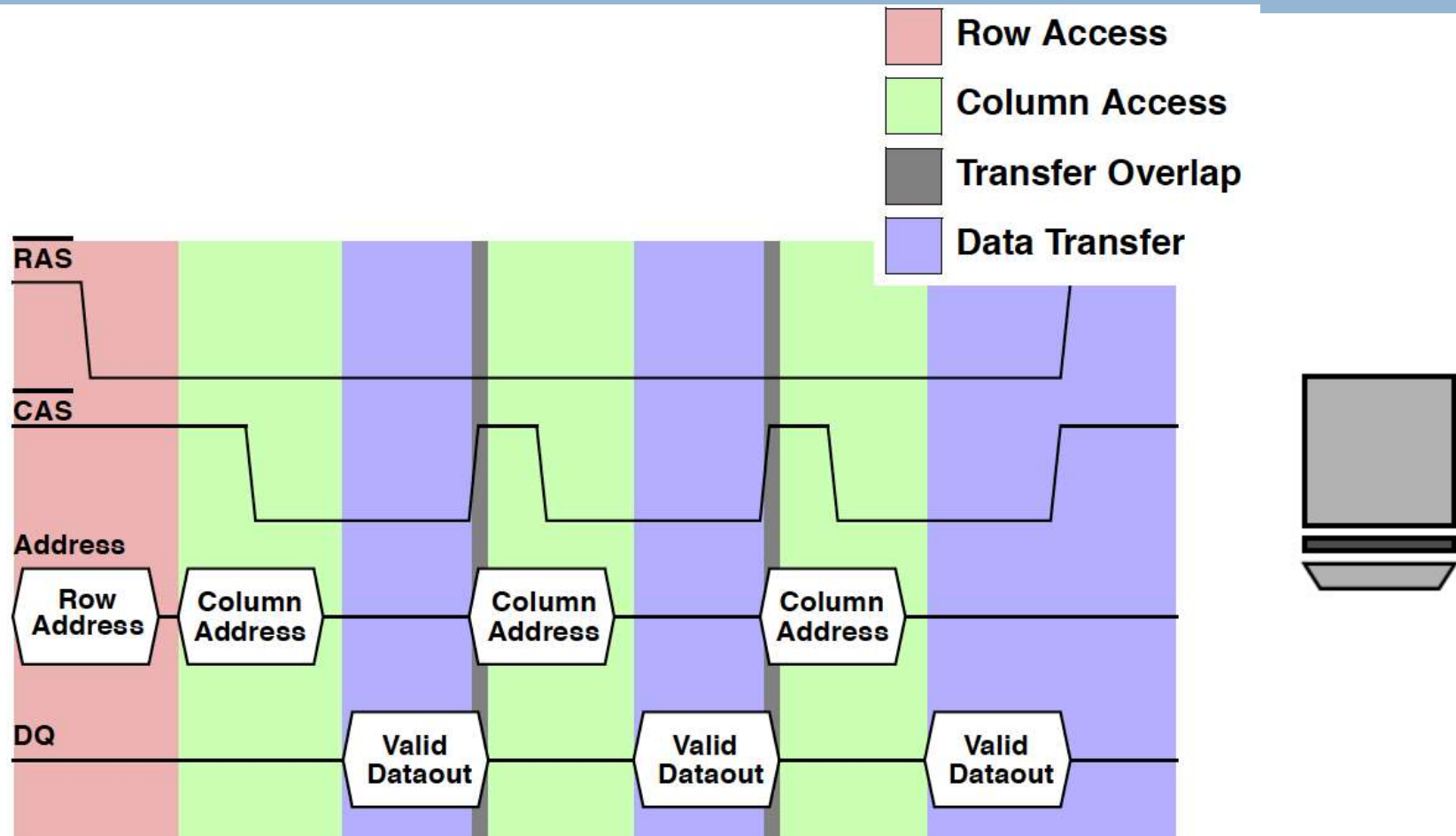
73



Source: “記憶體10年技術演進史，系統顆粒DDR與顯示顆粒GDDR差在哪？” by Tandee on internet

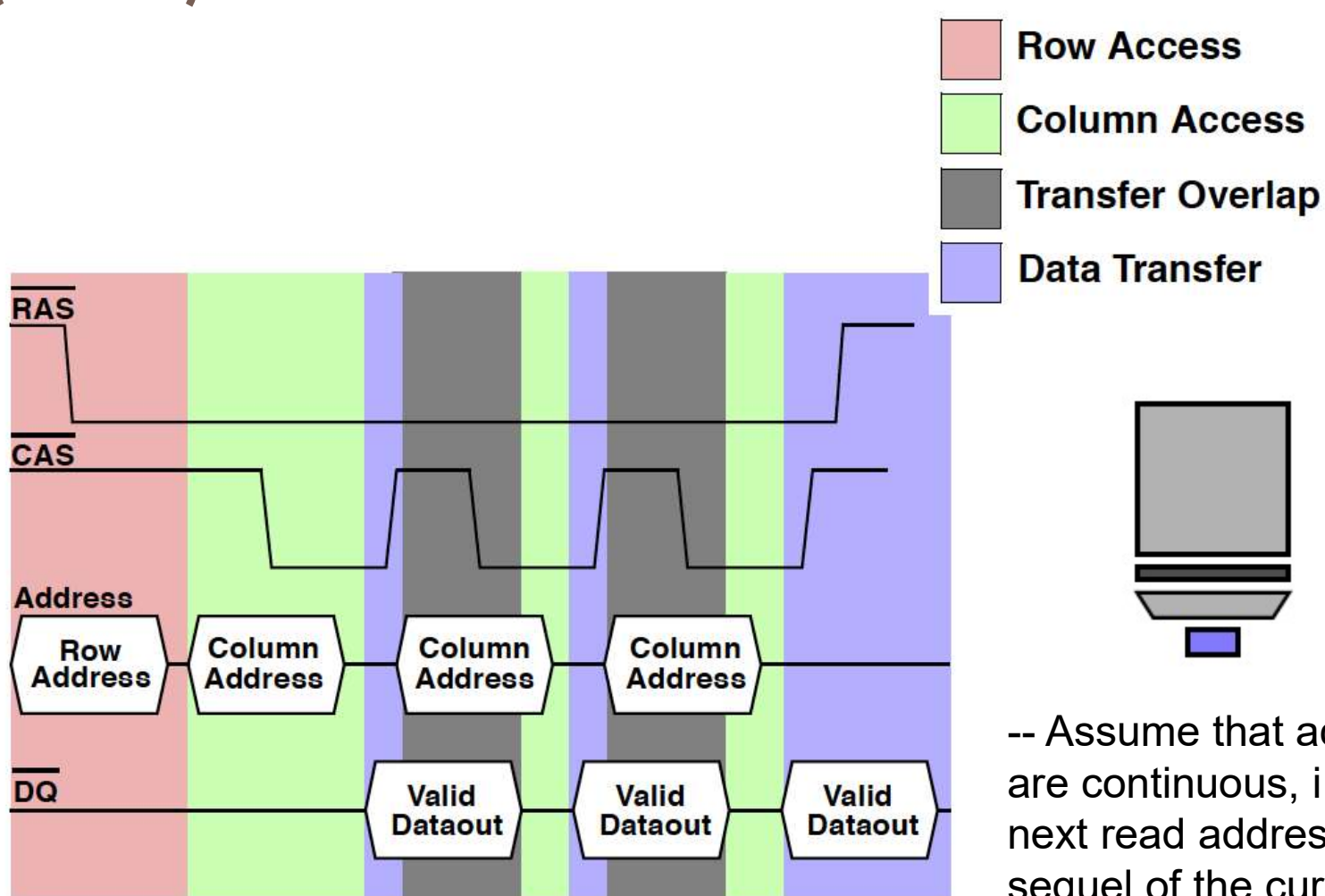
Reading Timing for Fast Page Mode

74



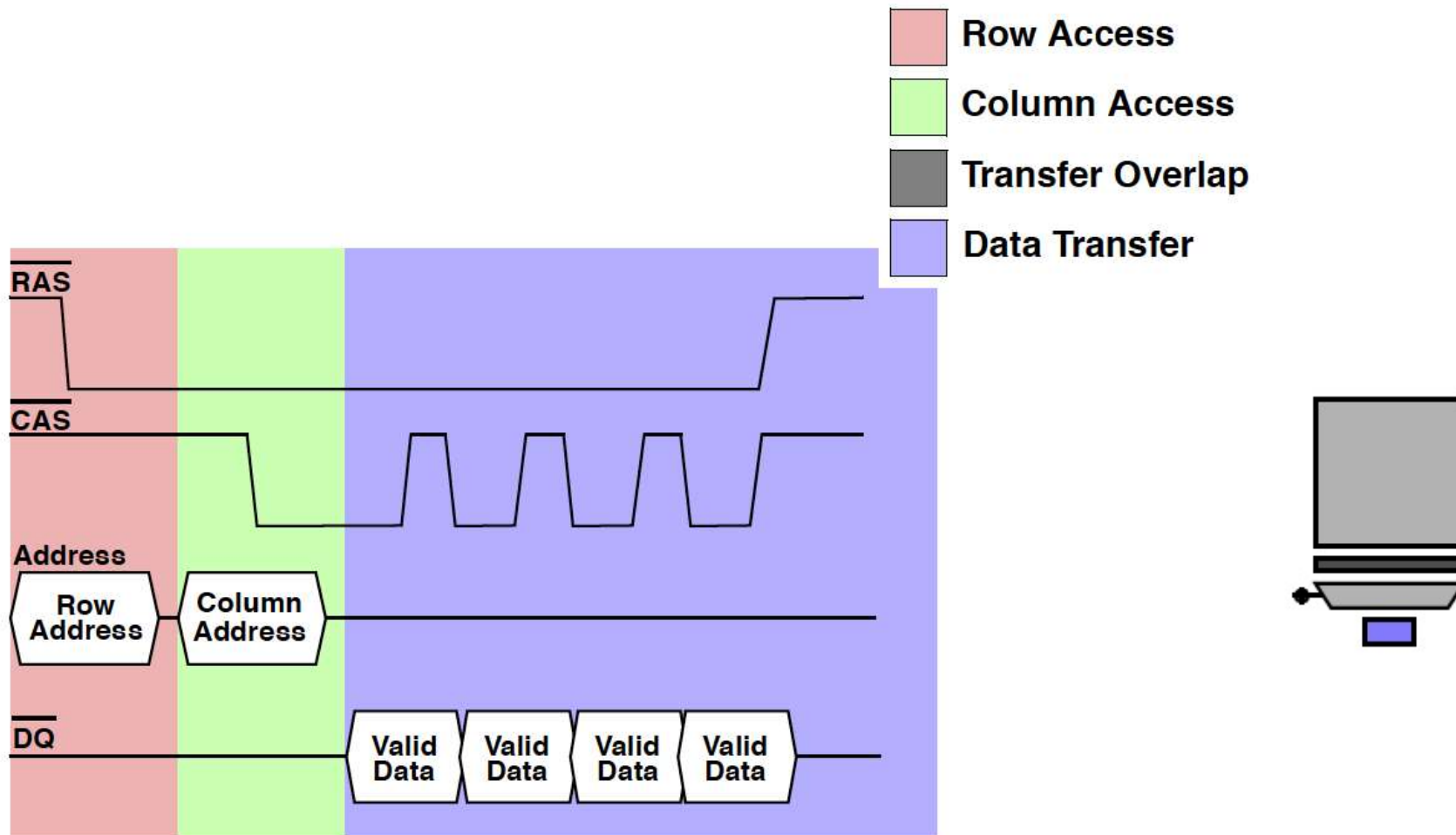
Read Timing for Extended Data Out (EDO)

75



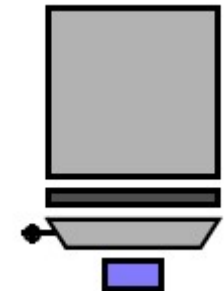
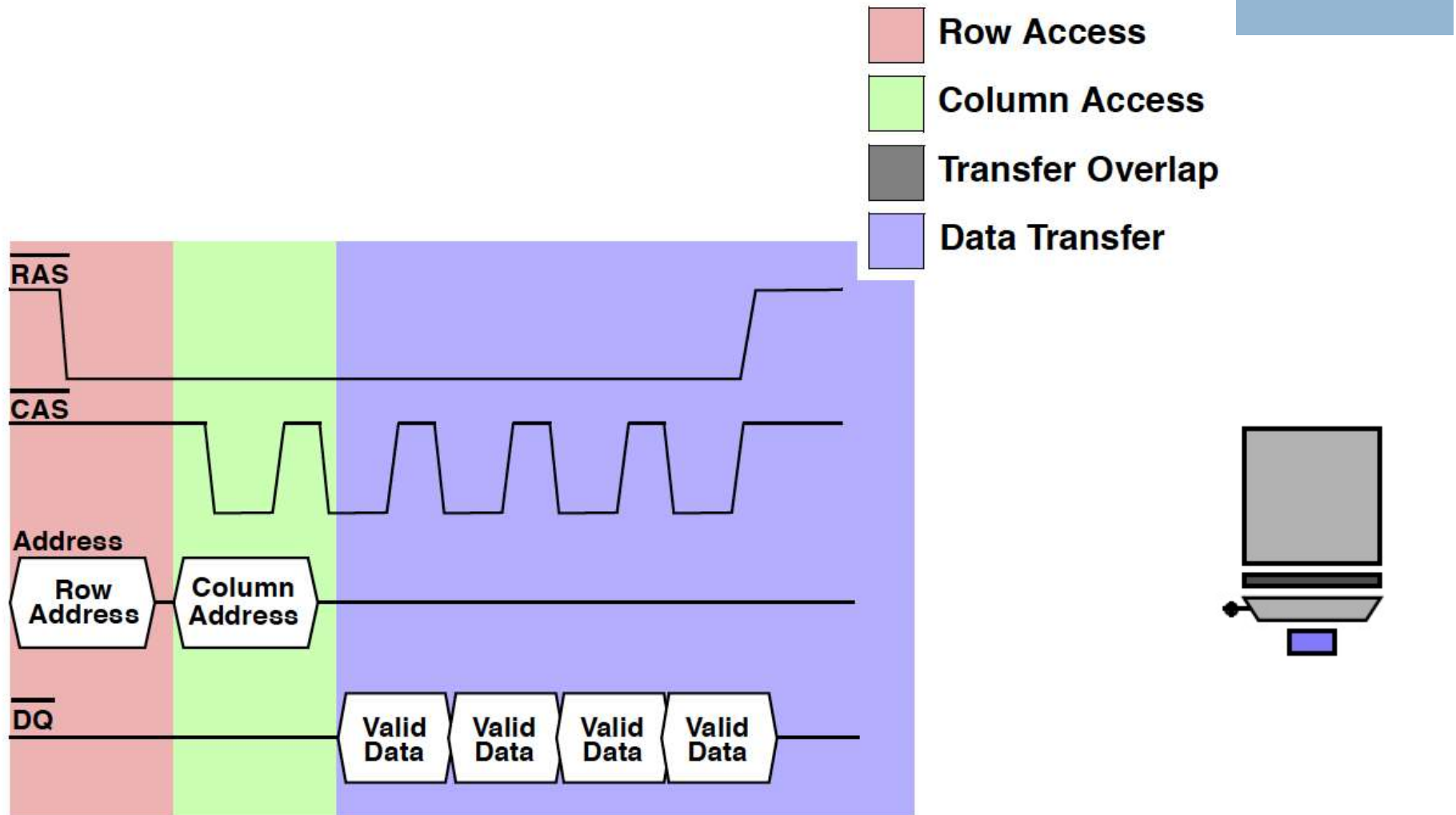
Read Timing for Burst EDO

76



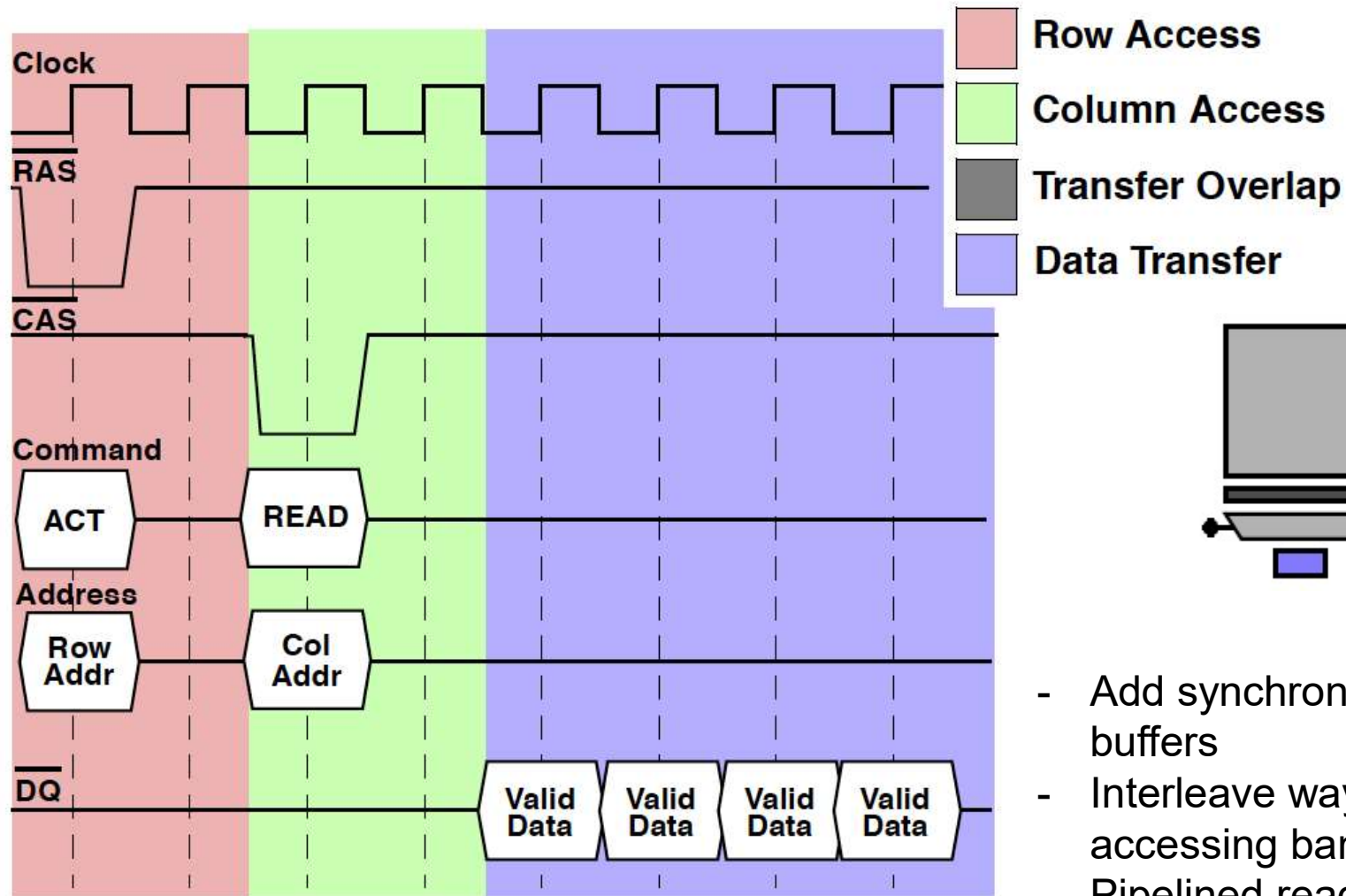
Read Timing for Pipeline Burst EDO

77



Read Timing for Synchronous DRAM

78



- Add synchronous dual buffers
- Interleave way in accessing banks
- Pipelined read or pipelined write

DDR SDRAM

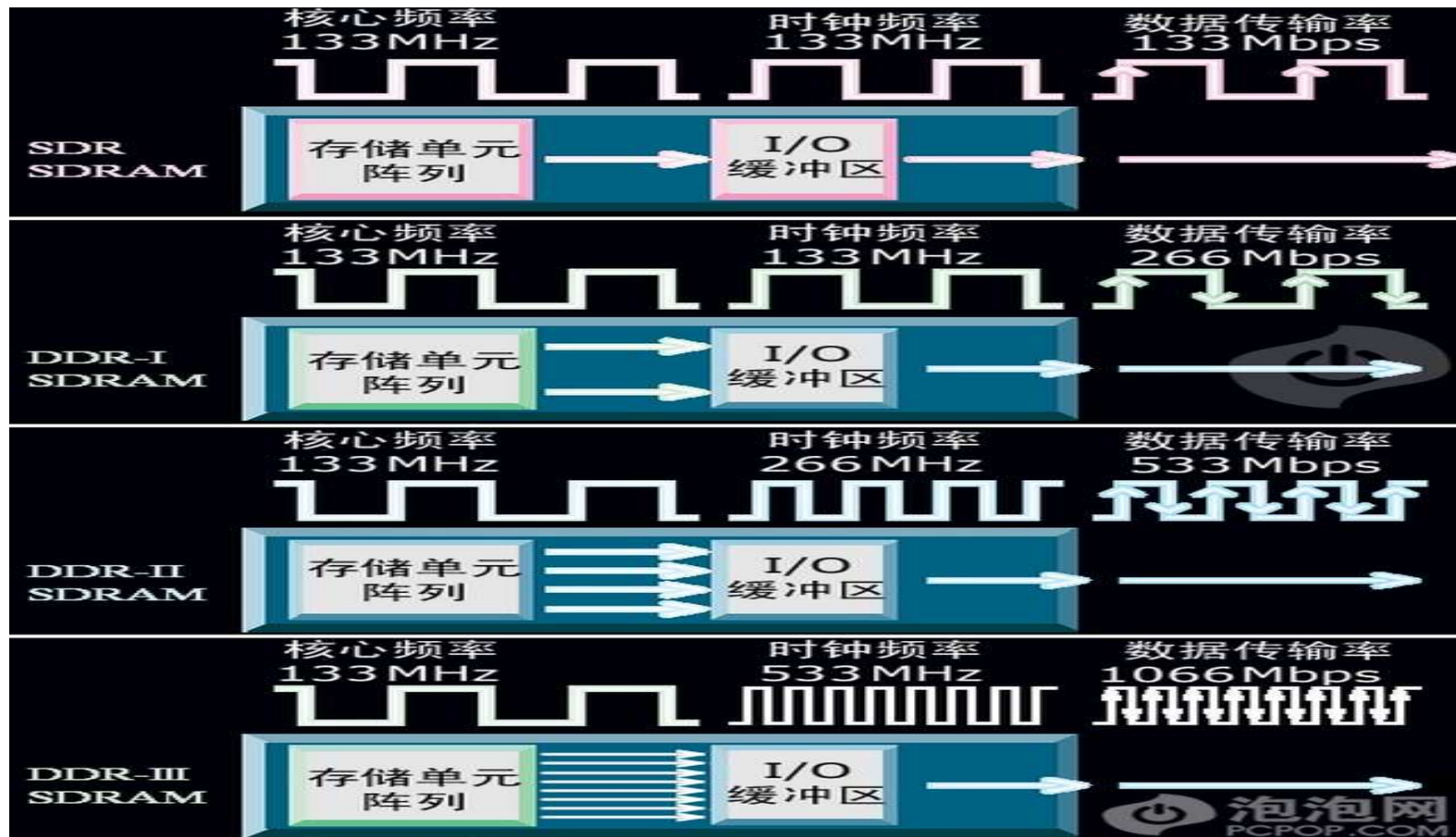
79

- Double Data Rate SDRAM
- Data can be transmitted on both clock edges

DDR SDRAM Standard	Bus clock (MHz)	Internal rate (MHz)	Prefetch (min burst)	Transfer Rate (MT/s)	Voltage	<u>DIMM</u> pins	<u>SO-DIMM</u> pins	<u>MicroDIMM</u> pins
DDR	100–200	100–200	2n	200–400	2.5/2.6	184	200	172
<u>DDR2</u>	200–533	100–266	4n	400–1066	1.8	240	200	214
<u>DDR3</u>	400–1066	100–266	8n	800–2133	1.5	240	204	214
<u>DDR4</u>	800–1200	200–300	8n	1600–2400	1.2	288	260	214

Clocking and Prefetch Buffering for Different Generations

80



Source: “記憶體10年技術演進史，系統顆粒DDR與顯示顆粒GDDR差在哪？” by Tandee on internet

GDDR (Graphics Double Data Rate)

81

版本	GDDR	GDDR2	gDDR2	GDDR3	gDDR3	GDDR4	GDDR5
預取量	2bit	4bit	4bit	4bit	8bit	8bit	8bit
對應世代	DDR	DDR2	DDR2	DDR2	DDR3	DDR3	DDR3
Burst	2/4/8bit	4/8bit	4/8bit	4/8bit	4/8bit	4/8bit	8bit
額定電壓	2.5V	2.5V	1.8V	1.8V	1.5V	1.5V	1.5V
單顆容量	16/32MB	32MB	32/64/128MB	32/64/128MB	64/128MB	64/128MB	64/128MB
介面頻寬	16/32bit	32bit	16bit	32bit	16bit	32bit	16/32bit
封裝針腳	66/144	144	84	136/144	96	136	170
Bank數量	2 / 4	4 / 8	4 / 8	4 / 8	8	8 / 16	8 / 16
等效時脈	300-900	800-1000	700-1200	1000-2600	1000-2000	2000-3000	3600-6000

Source: “記憶體10年技術演進史，系統顆粒DDR與顯示顆粒GDDR差在哪？” by Tandee on internet