

CS510 Computer Architecture

Lecture 18: DRAMs

Soontae Kim

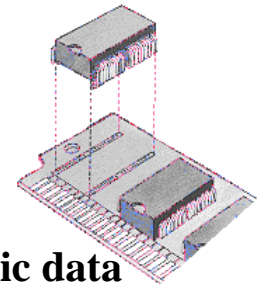
Spring 2017

School of Computing, KAIST

Announcements

- **Homework assignment #3**
 - Posted on class web site
 - Due on May 26 (Friday)
- **Term project**
 - Final presentations on June 7 and 9
 - Final report due on June 16

Main Memory



- Main memory generally utilizes Dynamic RAM (DRAM), which uses a single transistor and a capacitor to store a bit, but requires periodic data refreshes by reading every row. *capacitor discharge over time*
- Static RAM may be used for main memory if the added expense, low density, high power consumption, and complexity is feasible (e.g. Cray Vector Supercomputers).
- Main memory performance is affected by:
 - **Memory latency:** Affects cache miss penalty, M. Measured by:
 - **Access time:** The time it takes from the time a memory access request is issued to main memory to the time the requested word is available to cache/CPU.
 - **Cycle time:** The minimum time between requests to memory (greater than access time to allow address lines to be stable)
 - **Peak Memory bandwidth:** The maximum sustained data transfer rate between main memory and cache/CPU.
 - Realistic memory bandwidth < peak memory bandwidth



Logical DRAM Chip Organization

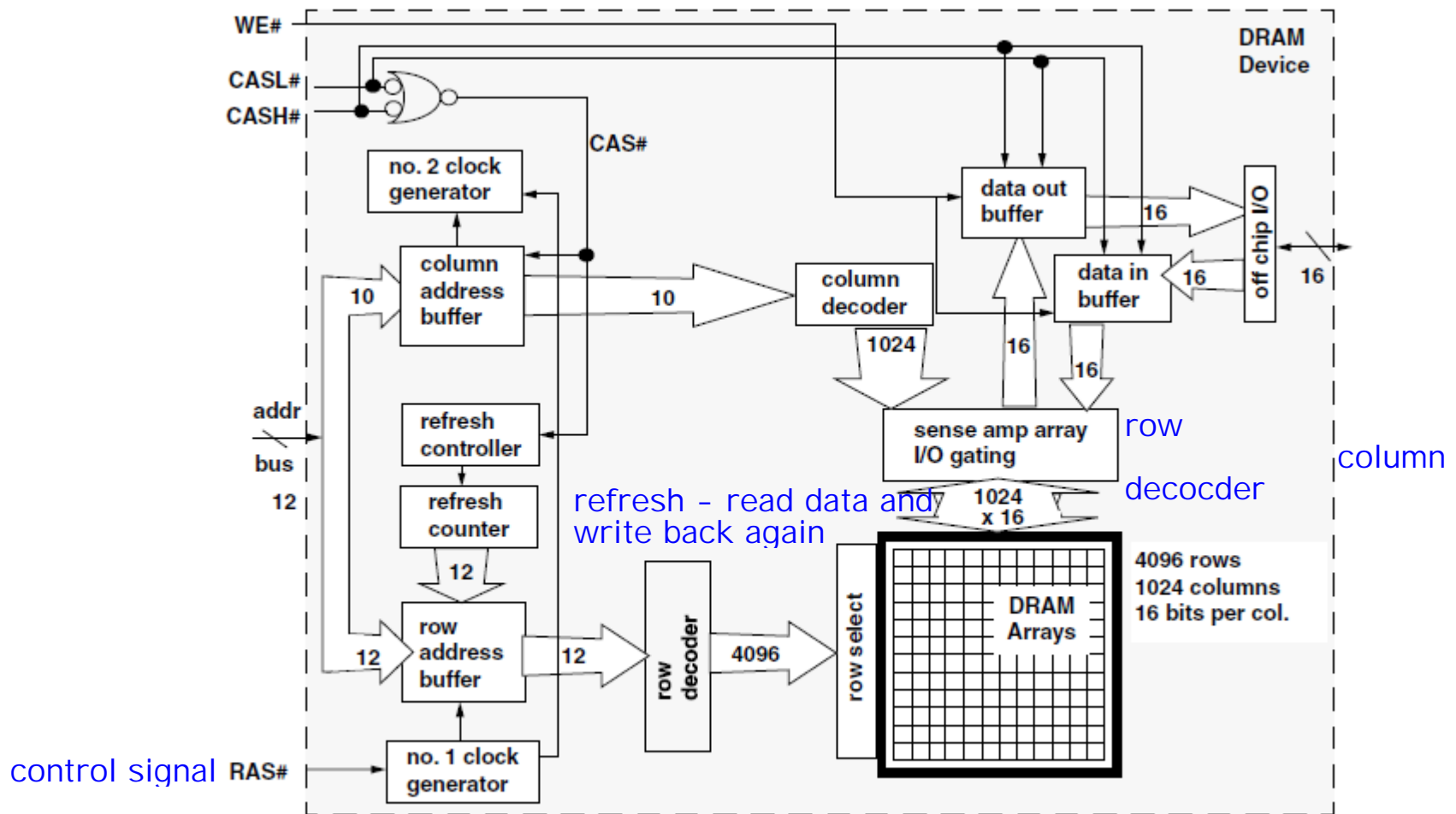


FIGURE 8.1: A 64-Mbit Fast Page Mode DRAM device (4096 x 1024 x 16).

Control Signals:

- 1 - Row Access Strobe (RAS): Low to latch row address
- 2 - Column Access Strobe (CAS): Low to latch column address

A periodic data refresh is required by reading every bit

Four Key DRAM Timing Parameters

- **t_{RAC}** : Minimum time from RAS (Row Access Strobe) line falling (activated) to the valid data output.
random access

- Used to be quoted as the nominal speed of a DRAM chip
- For a typical 64Mb DRAM $t_{\text{RAC}} = 60 \text{ ns}$

- **t_{RC}** : Minimum time from the start of one row access to the start of the next (memory cycle time).
random access cycle time

- $t_{\text{RC}} = t_{\text{RAC}} + \text{RAS Precharge Time}$ (bitline precharging time & bus reset electronic issue)
- $t_{\text{RC}} = 110 \text{ ns}$ for a 64Mbit DRAM with a t_{RAC} of 60 ns

- **t_{CAC}** : Minimum time from CAS (Column Access Strobe) line falling to valid data output.
column or page access = CAS latency tcl

- 12 ns for a 64Mbit DRAM with a t_{RAC} of 60 ns

- **t_{PC}** : Minimum time from the start of one column access to the start of the next.
Fast Page Cycle Time

- $t_{\text{PC}} = t_{\text{CAC}} + \text{CAS Precharge Time}$
- About 25 ns for a 64Mbit DRAM with a t_{RAC} of 60 ns

Simplified DRAM Speed Parameters

- Row Access Strobe (RAS) Time: (similar to t_{RAC}):
 - Minimum time from RAS (Row Access Strobe) line falling (activated) to the first valid data output.
 - A major component of memory latency.
 - Only improves ~ 5% every year.
- Column Access Strobe (CAS) Time/data transfer time: (similar to t_{CAC})
 - The minimum time required to read additional data by changing column address while keeping the same row address.
 - Along with memory bus width, determines peak memory bandwidth.

t_{RAS} : minimum time from RAS line falling to the first valid data output

t_{CAS} : minimum time from CAS line falling to the first valid data output

DRAM Generations

Year	Size	RAS (ns)	CAS (ns)	Cycle Time	Memory Type	
1980	64 Kb	150-180	75	250 ns	Page Mode	Asynchronous DRAM
1983	256 Kb	120-150	50	220 ns	Page Mode	
1986	1 Mb	100-120	25	190 ns		
1989	4 Mb	80-100	20	165 ns	Fast Page Mode	
1992	16 Mb	60-80	15	120 ns	EDO	
1996	64 Mb	50-70	12	110 ns	PC66 SDRAM	Synchronous DRAM
1998	128 Mb	50-70	10	100 ns	PC100 SDRAM	
2000	256 Mb	45-65	7	90 ns	PC133 SDRAM	
2002	512 Mb	40-60	5	80 ns	PC2700 DDR	
2004	1Gb	35-55	5	70 ns	PC4300 DDR2	
2006	2Gb	30-50	2.5	60 ns	PC8500 DDR3	
	32000:1 (Capacity)		30:1 (~bandwidth)	4:1 (Latency)	frequency bandwidth	

Basic Memory Bandwidth Improvement/Miss Penalty (M)

Reduction Techniques

- **Wider Main Memory (CPU-Memory Bus):**

Memory bus width is increased to a number of words (usually up to the size of a cache block).

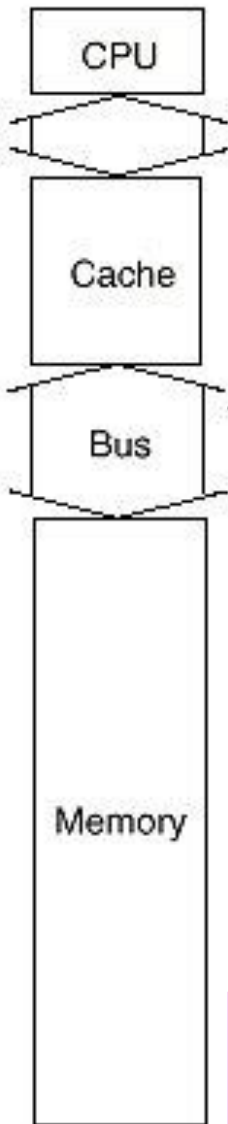
- Memory bandwidth is proportional to memory bus width.
 - e.g Doubling the width of cache and memory doubles potential memory bandwidth available to the CPU.
- The miss penalty is reduced since fewer memory bus accesses are needed to fill a cache block on a miss.

- **Interleaved (Multi-Bank) Memory:**

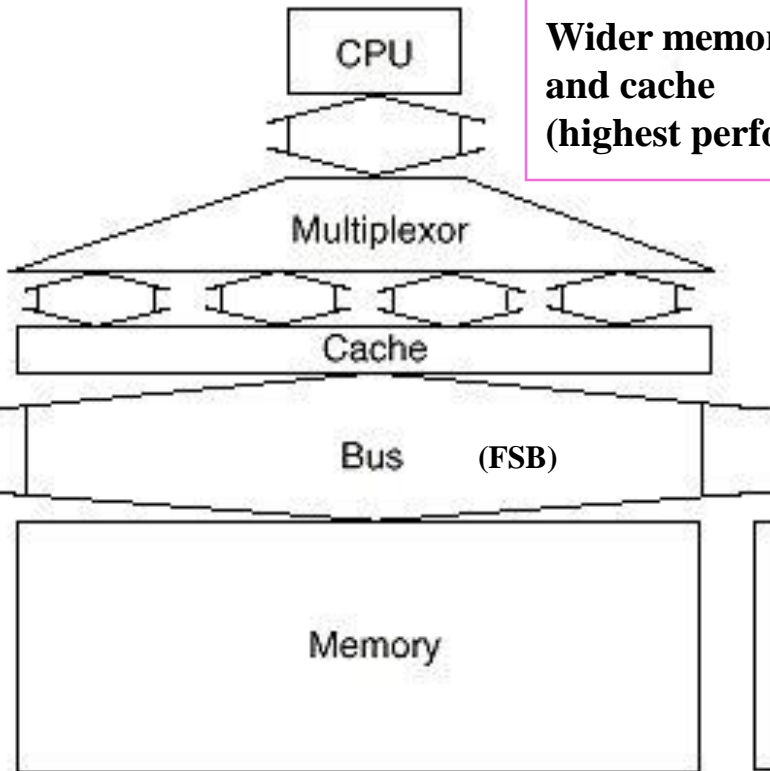
Memory is organized as a number of independent banks.

- Multiple interleaved memory reads or writes are accomplished by sending memory addresses to several memory banks at once or pipeline accesses to the banks.
- **Interleaving factor:** Refers to the mapping of memory addressees to memory banks. Goal reduce bank conflicts.
e.g. using 4 banks (one word wide), bank 0 has all words whose addresses are: $(\text{word address mod } 4) = 0$

(a) One-word-wide memory organization

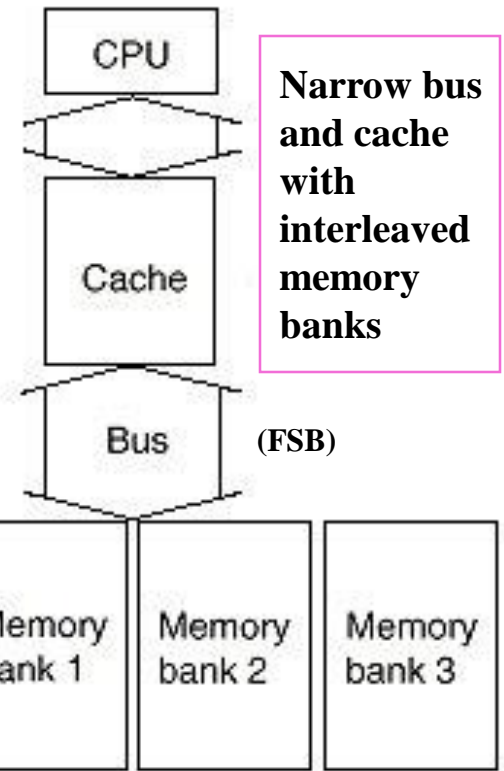


(b) Wide memory organization



**Wider memory, bus and cache
(highest performance)**

(c) Interleaved memory organization



Narrow bus and cache with interleaved memory banks

Three examples of bus width, memory width, and memory interleaving to achieve higher memory bandwidth

**Simplest design:
Everything is the width of one word (lowest performance)**

Front Side Bus (FSB) = System Bus = CPU-memory Bus

Four Way (Four Banks) Interleaved Memory

Memory Bank Number

	Bank 0	Bank 1	Bank 2	Bank 3
Address	0	1	2	3
Within	4	5	6	7
Bank	8	9	10	11
	12	13	14	15
	16	17	18	19
	20	21	22	23

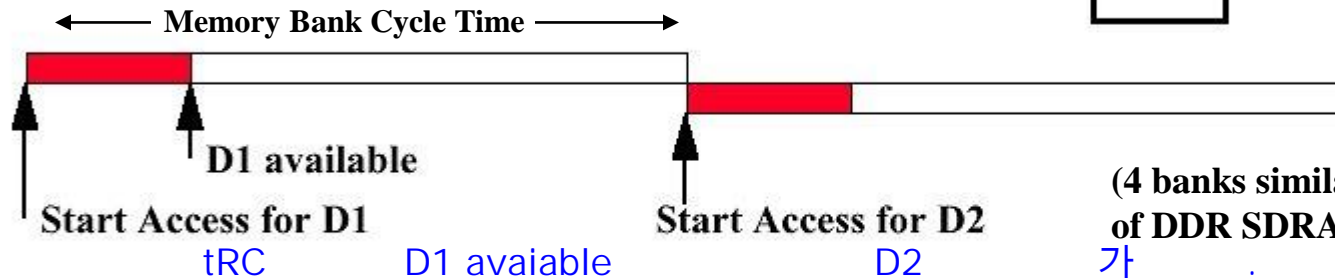
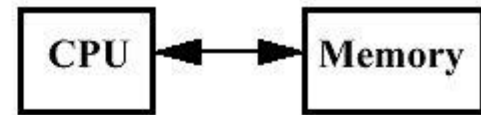
Bank Width = One Word

Bank Number = (Word Address) Mod (4)

Memory Bank Interleaving

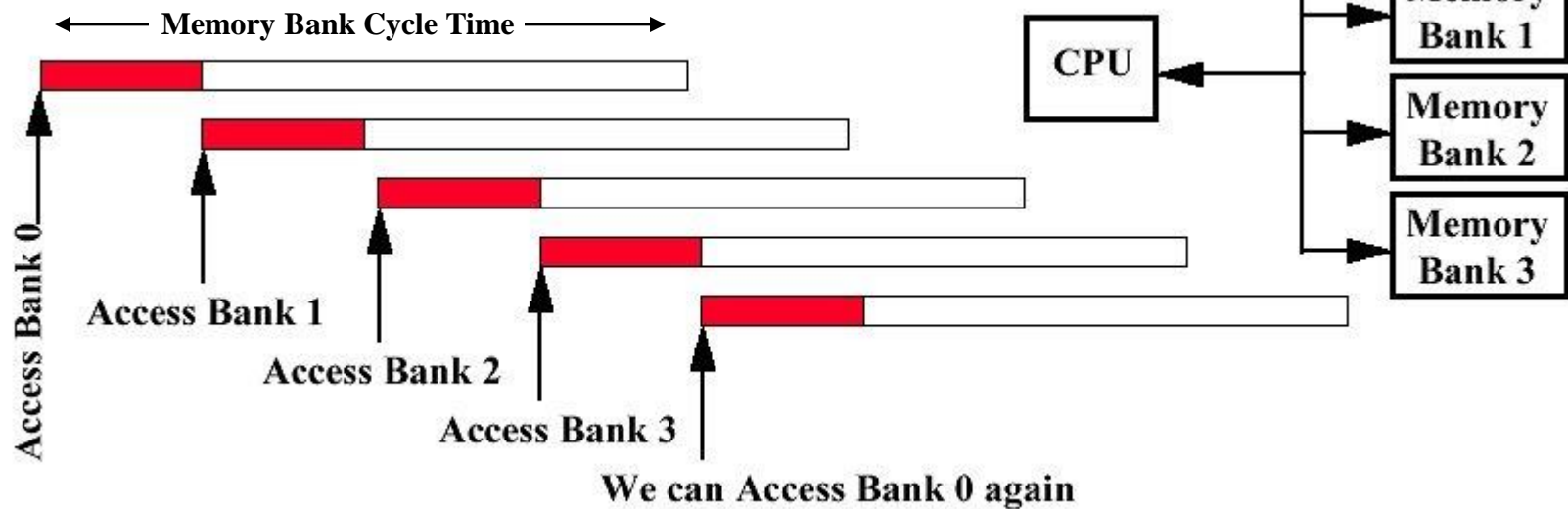
Access Pattern without Interleaving:

(One Bank)



Pipeline access to different memory banks to increase effective bandwidth

Access Pattern with 4-way Interleaving:



Bank interleaving does not reduce latency of accesses to the same bank

Synchronous Dynamic RAM, (SDRAM) (mid 90s)

Organization

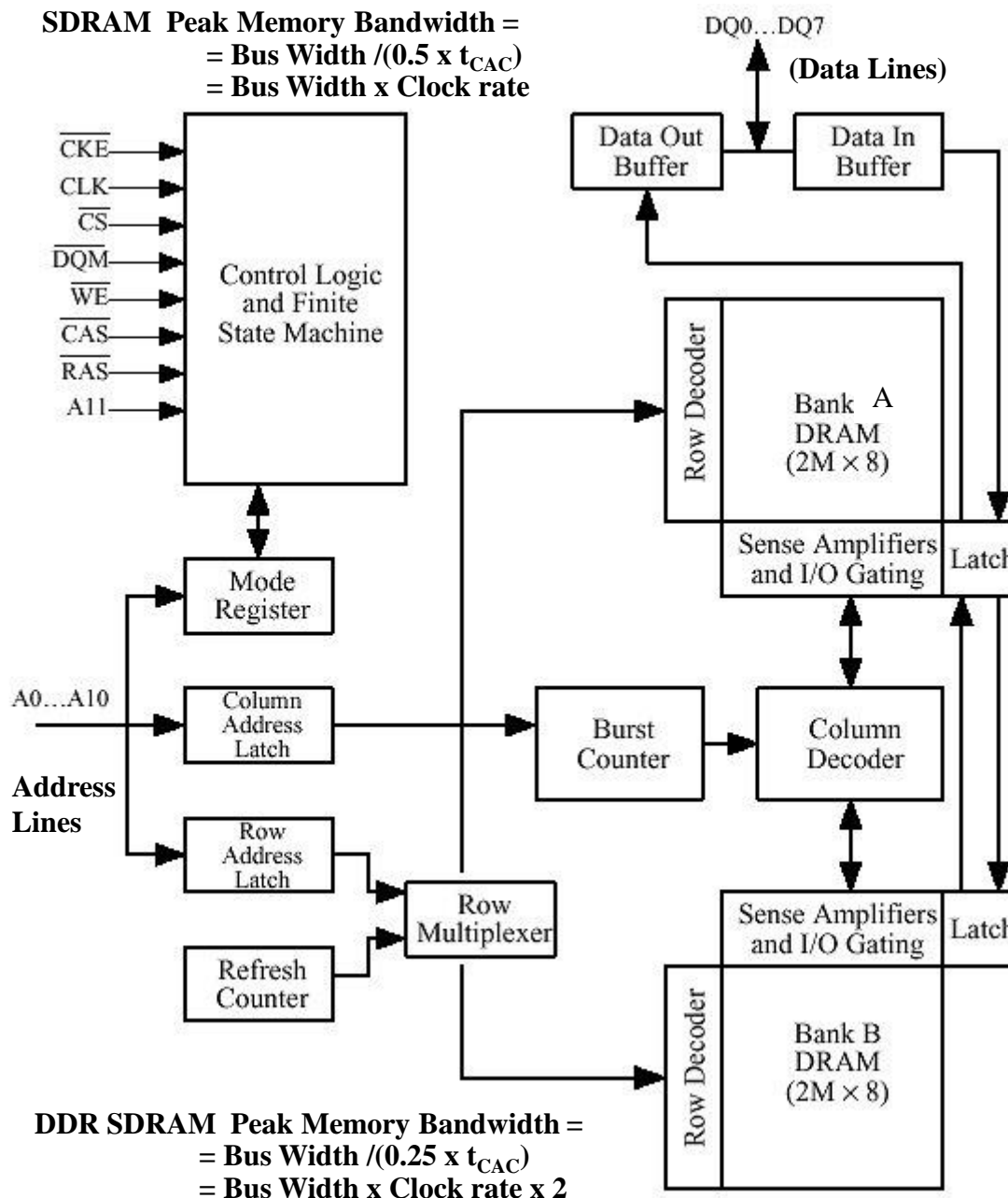
SDRAM speed is rated at max. clock speed supported:
 100MHZ = PC100
 133MHZ = PC133

DDR SDRAM (late 90s - current)

organization is similar but **four banks** are used in each DDR SDRAM chip instead of two.

Data transfer on both **rising and falling edges of the clock**

DDR SDRAM rated by maximum memory bandwidth if the bus width = 8 byte
 PC3200 = 8 bytes x 200 MHz x 2
 = 3200 Mbytes/sec



DDR SDRAM Standard	Release year	Bus clock (MHz)	Internal rate (MHz)
DDR1	2000	100–200	100–200
DDR2	2003	200–533.33	100–266.67
DDR3	2007	400–1066.67	100–266.67
DDR4	2014	1066.67–2133.33	133.33–266.67

DDR SDRAM

the prefetch buffer depth can also be thought of as the ratio between the core memory frequency (internal rate) and IO frequency (bus clock).
In DDR3, the IOs will operate 8 times faster than the memory core

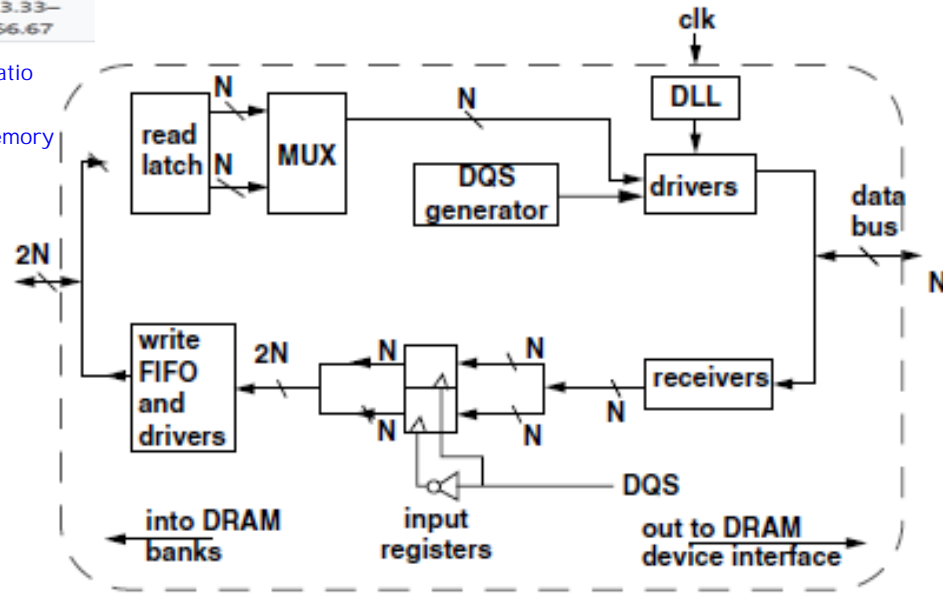


FIGURE 12.17: DDR SDRAM device I/O.

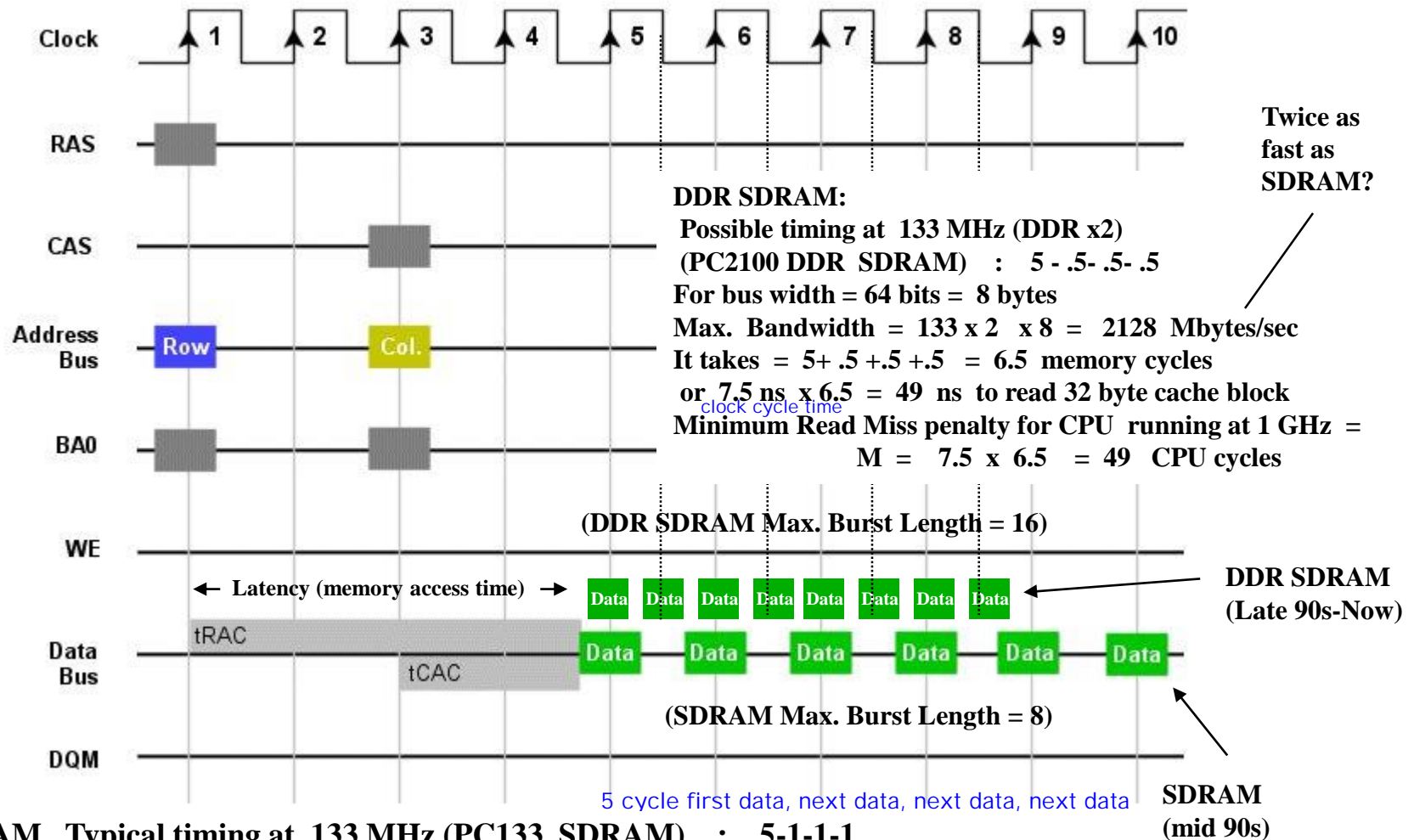
- **M-bit prefetch architecture** simplifies the subsequent column access by allowing a single address request to result in multiple data words
 - Data rate multiplication architecture
 - M represents multiplication factor b/t DRAM device's internal width of data movement and width of the data bus on the device interface
 - DDR (DDR II, DDR III) SDRAM supports 2-bit (4-bit, 8-bit) prefetch architecture

In a prefetch buffer architecture, when a memory access occurs to a row the buffer grabs a set of adjacent data words on the row and reads them out in rapid-fir sequence on the IO pins, without the need for individual column address requests.

DDR1 : two adjacent data words will be read from each chip in the same clock cycle and placed in the prefetch buffer. each word will then be transmitted on consecutive rising and falling edges of the clock cycle.

DDR2 with a 4n prefetch buffer : four consecutive data words are read and placed in buffer while a clock, which is twice faster than the external clock of DDR, transmits each of the word in consecutive rising and falling edge of the faster external clock.

SDRAM Read Simplified SDRAM/DDR SDRAM Read Timing



SDRAM Typical timing at 133 MHz (PC133 SDRAM) : 5-1-1-1

For bus width = 64 bits = 8 bytes Max. Bandwidth = 133 x 8 = 1064 Mbytes/sec

It takes = 5+1+1+1 = 8 memory cycles or 7.5 ns x 8 = 60 ns to read 32 byte cache block

Minimum Read Miss penalty for CPU running at 1 GHz = M = 7.5 x 8 = 60 CPU cycles

In this example for SDRAM: M = 60 cycles for DDR SDRAM: M = 49 cycles

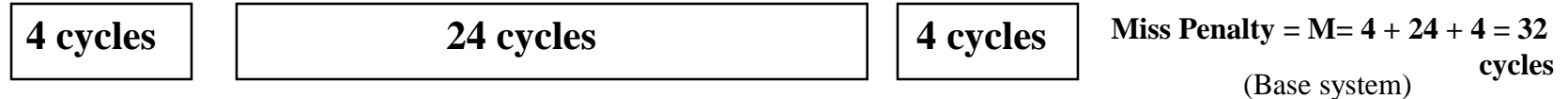
Thus accounting for access latency DDR is 60/49 = 1.22 times faster

Not twice as fast (2128/1064 = 2) as indicated by peak bandwidth!

Memory Width, Interleaving: Performance Example

Given the following system parameters with single unified cache level L_1 (ignoring write policy):

Block size= 1 word Memory bus width= 1 word Miss rate =3% M = Miss penalty = 32 cycles
(4 cycles to send address 24 cycles access time, 4 cycles to send a word to CPU)



Memory access/instruction = 1.2 $CPI_{\text{execution}}$ (ignoring cache misses) = 2

Miss rate (block size = 2 word = 8 bytes) = 2% Miss rate (block size = 4 words = 16 bytes) = 1%

- The CPI of the base machine with 1-word blocks = $2 + (1.2 \times 0.03 \times 32) = 3.15$

Increasing the block size to two words (64 bits) gives the following CPI: (miss rate = 2%)

- 32-bit bus and memory, no interleaving, $M = 2 \times 32 = 64$ cycles $CPI = 2 + (1.2 \times .02 \times 64) = 3.54$
- 32-bit bus and memory, interleaved, $M = 4 + 24 + 8 = 36$ cycles $CPI = 2 + (1.2 \times .02 \times 36) = 2.86$
- 64-bit bus and memory, no interleaving, $M = 32$ cycles $CPI = 2 + (1.2 \times 0.02 \times 32) = 2.77$

Increasing the block size to four words (128 bits); resulting CPI: (miss rate = 1%)

- 32-bit bus and memory, no interleaving , $M = 4 \times 32 = 128$ cycles $CPI = 2 + (1.2 \times 0.01 \times 128) = 3.54$
- 32-bit bus and memory, interleaved , $M = 4 + 24 + 16 = 44$ cycles $CPI = 2 + (1.2 \times 0.01 \times 44) = 2.53$
- 64-bit bus and memory, no interleaving, $M = 2 \times 32 = 64$ cycles $CPI = 2 + (1.2 \times 0.01 \times 64) = 2.77$
- 64-bit bus and memory, interleaved, $M = 4 + 24 + 8 = 36$ cycles $CPI = 2 + (1.2 \times 0.01 \times 36) = 2.43$
- 128-bit bus and memory, no interleaving, $M = 32$ cycles $CPI = 2 + (1.2 \times 0.01 \times 32) = 2.38$

Miss Penalty = M = Number of CPU stall cycles for an access missed in cache and satisfied by main memory