

Memory Systems

Memory Hierarchy

- Smaller and closer are faster
 - Flipflops
 - Register files
- Larger and further away are cheaper
 - DRAM
 - Hard disks

STORAGE	DELAYS: Absolute	DELAY: in CPU clocks
Fastest	Isolated latches	a few psecs to a few nsecs
	Registers/ Small RAM	a few 10s of psecs to a few nsecs.
	Large on-chip SRAM	a few nsecs.
Off chip SRAM	a few nsecs. to 20 nsecs.	
Off-chip DRAM	40 nsecs. and higher	NEW non-volatile: P-RAM, STT-RAM*
Flash storage	tens of microsecs, slower writes	
Magnetic disks	6 to 14 msecs.	
Optical disks	10 to 20 msecs.	
Slowest		
 :on-chip storage resources		

* price point unclear

Examples:

- ▶ Register file bandwidth for a 4-issue per cycle superscalar CPU running at 3 GHz.:
 - 8 operands, 32-bits wide per cycle => $BW = 96 \text{ Gbits/sec.}$
 - Double this for 64-bit datapaths
- ▶ CPU I/O pin bandwidth: 128 I/O pins, 1 GHz. external memory bus: 32 Gbits/sec. – this is the best data rate one can realize from off-chip memory devices. Can double this by using two transfers per clock (one per clock edge, as in DDR memory interfaces)
- ▶ Typical disk transfer rate: several 10s to 125 MBytes/sec. (single disk unit, sustained data transfer rate)

- ▶ DRAM bandwidth: regular DRAM chips –
 - Single byte access (“byte-wide” DRAM chip): 1 byte/60 nsecs.
 - Internal BW based on a 1024-bit row size: 1K to 4K bits/40 nsecs.
 - Significant loss in data rate occurs as only a byte gets selected internally within the chip and gets sent out!
 - Burst or page mode data rate from DRAM is faster (roughly about 2 times individual byte access rates)

Characteristics of Memory Traffic

- Memory instructions vs. all instructions:
 - RISC: 15% to 20%
 - CISC (x86): 35%
- Loads vs. Stores: 2 to 1
- Caches make external memory access bursty (smoother would be better)
- Strided accesses in many applications

Memory Performance

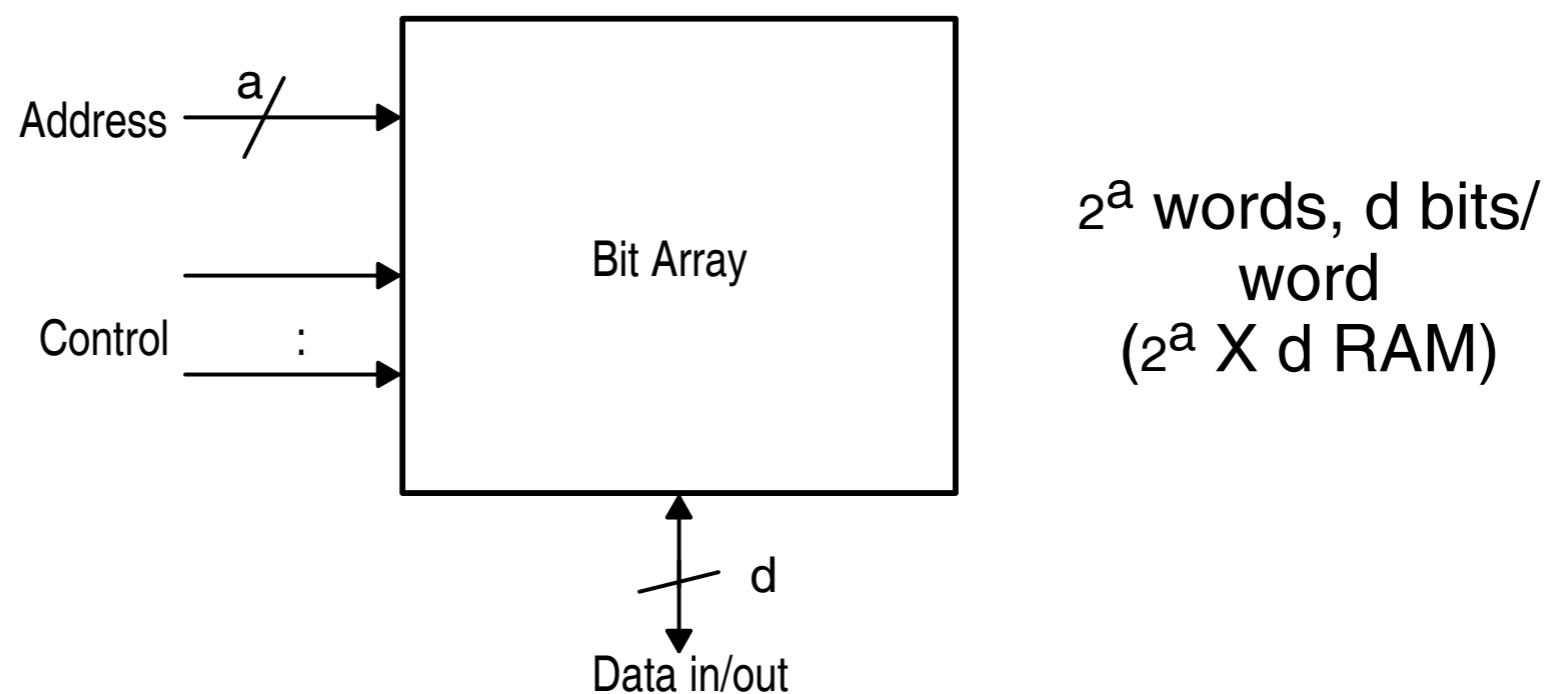
- Processor-memory performance gap
- Latency = Request & Response time
- Bandwidth = Steady-state transfer speed
- Effective transfer time:
 $\text{latency} + (\text{data_size} / \text{bandwidth})$
- Bandwidth easier to improve than latency

Techniques for Improving Memory Performance

- Memory interleaving
- Cache memory
- Newer memory interfaces
- Prefetching
- Stream buffers
- Store bypassing by loads and predicting store bypassing (“dynamic memory disambiguation”)
- Cache miss prediction
- Simultaneous multithreading

RAM Devices

- Basic storage component: *bitcell* – one bit of storage:
 - ▶ Static RAM (SRAM) bitcell = flip-flop (back-to-back connected inverters)
 - ▶ Dynamic (DRAM) bitcell = MOS capacitor: charged => 1, discharged => 0
- The RAM as a black box:

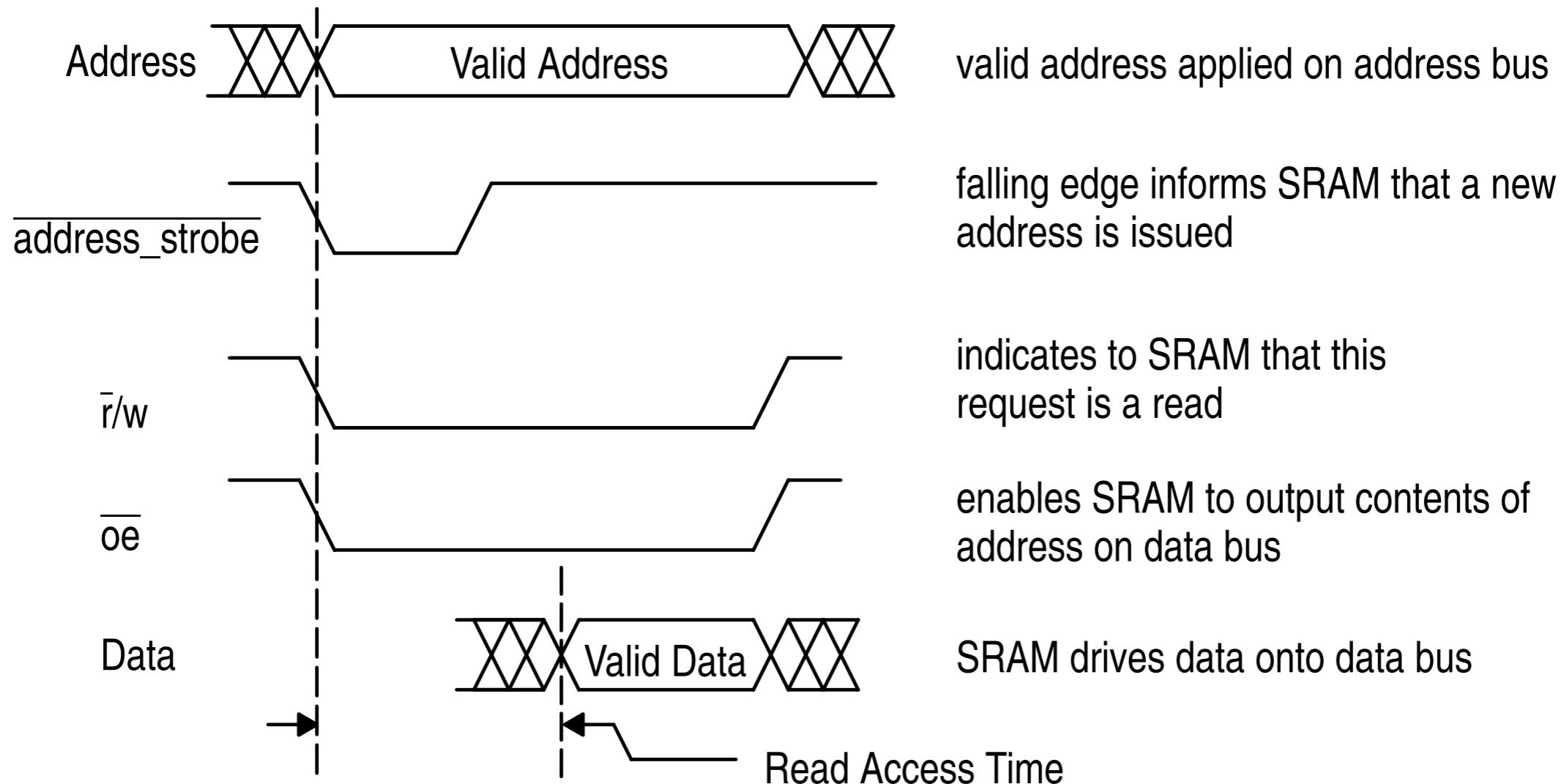




CMOS Static RAM
1 Meg (64K x 16-Bit)

A4	1		44	A5
A3	2		43	A6
A2	3		42	A7
A1	4		41	OE
A0	5		40	BHE
CS	6		39	BLE
I/O 0	7		38	I/O 15
I/O 1	8		37	I/O 14
I/O 2	9		36	I/O 13
I/O 3	10		35	I/O 12
Vcc	11	SO44-1	34	Vss
Vss	12	SO44-2	33	Vcc
I/O 4	13		32	I/O 11
I/O 5	14		31	I/O 10
I/O 6	15		30	I/O 9
I/O 7	16		29	I/O 8
WE	17		28	NC
A15	18		27	A8
A14	19		26	A9
A13	20		25	A10
A12	21		24	A11
NC	22		23	NC

Asynchronous static RAM chip

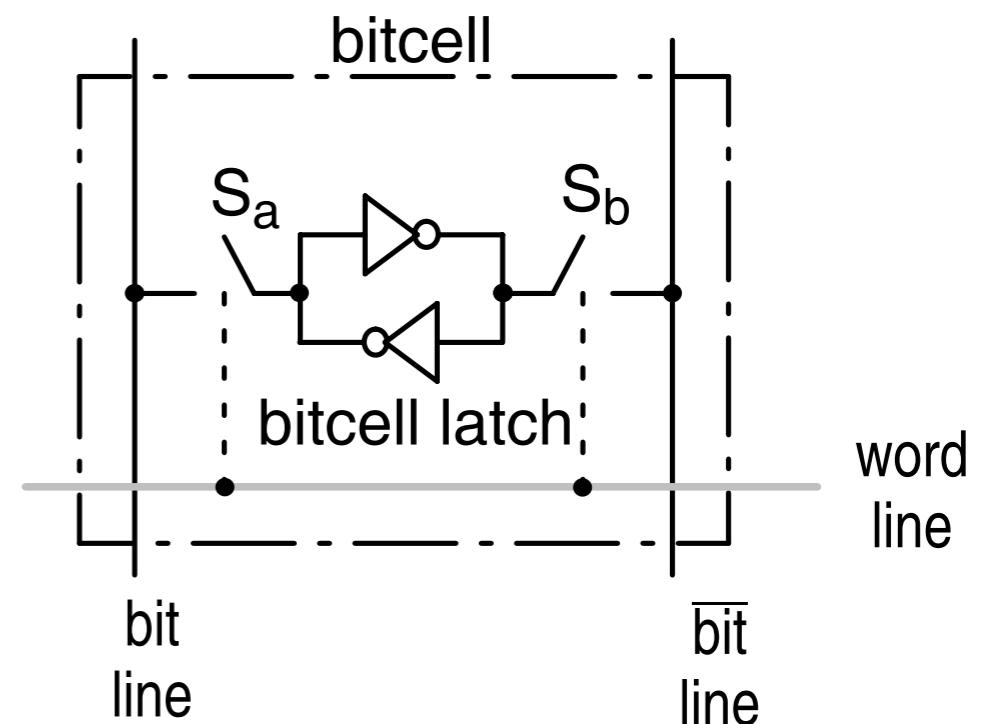
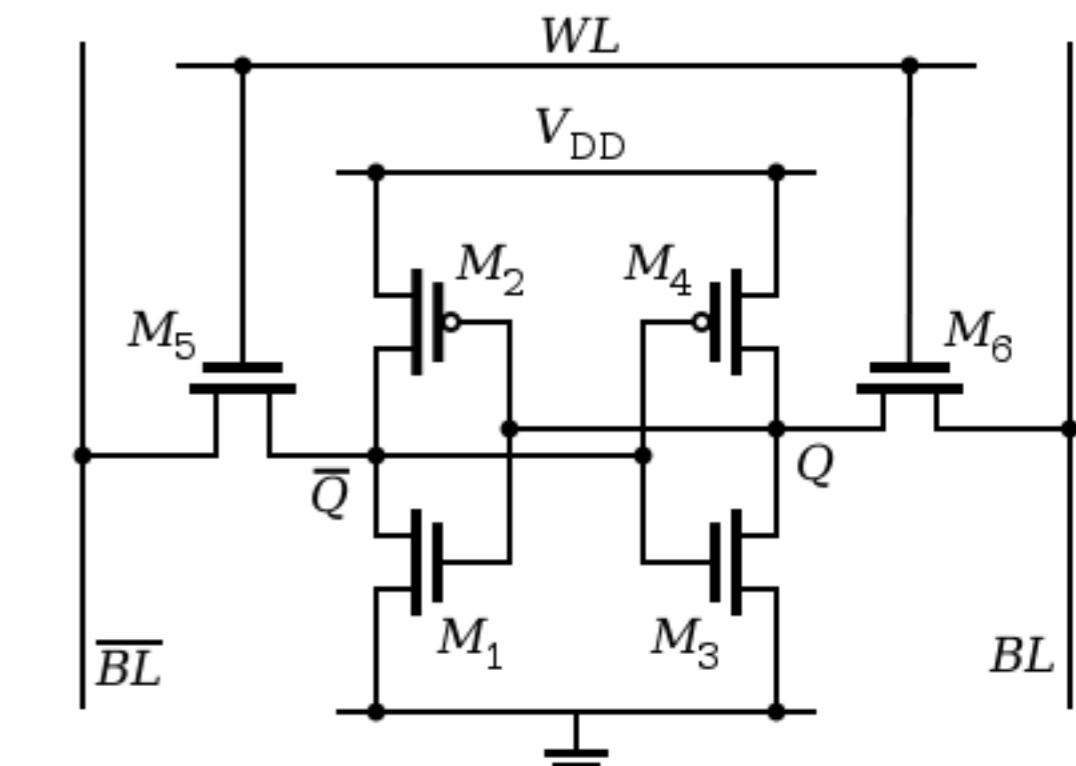
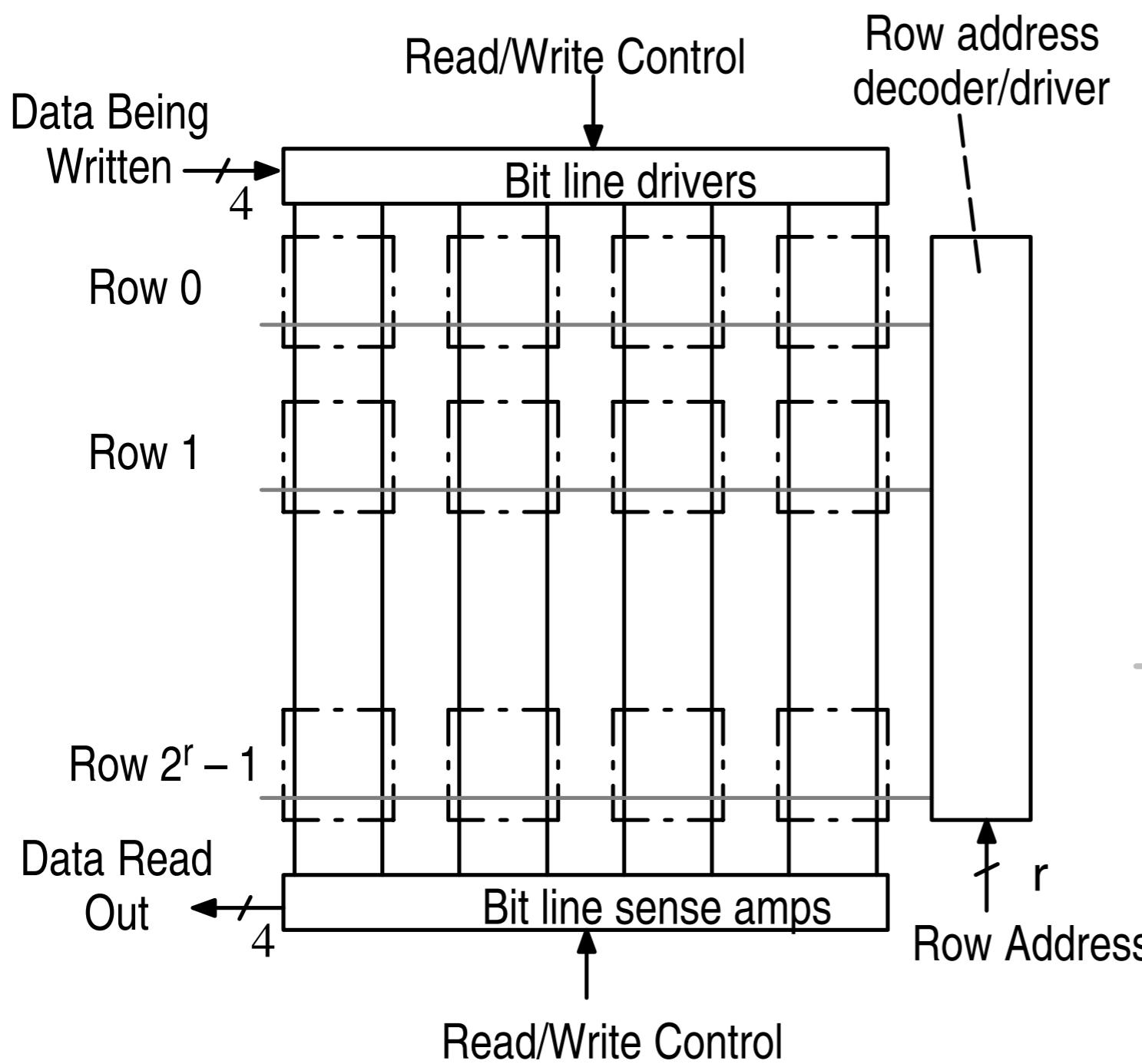


= undefined

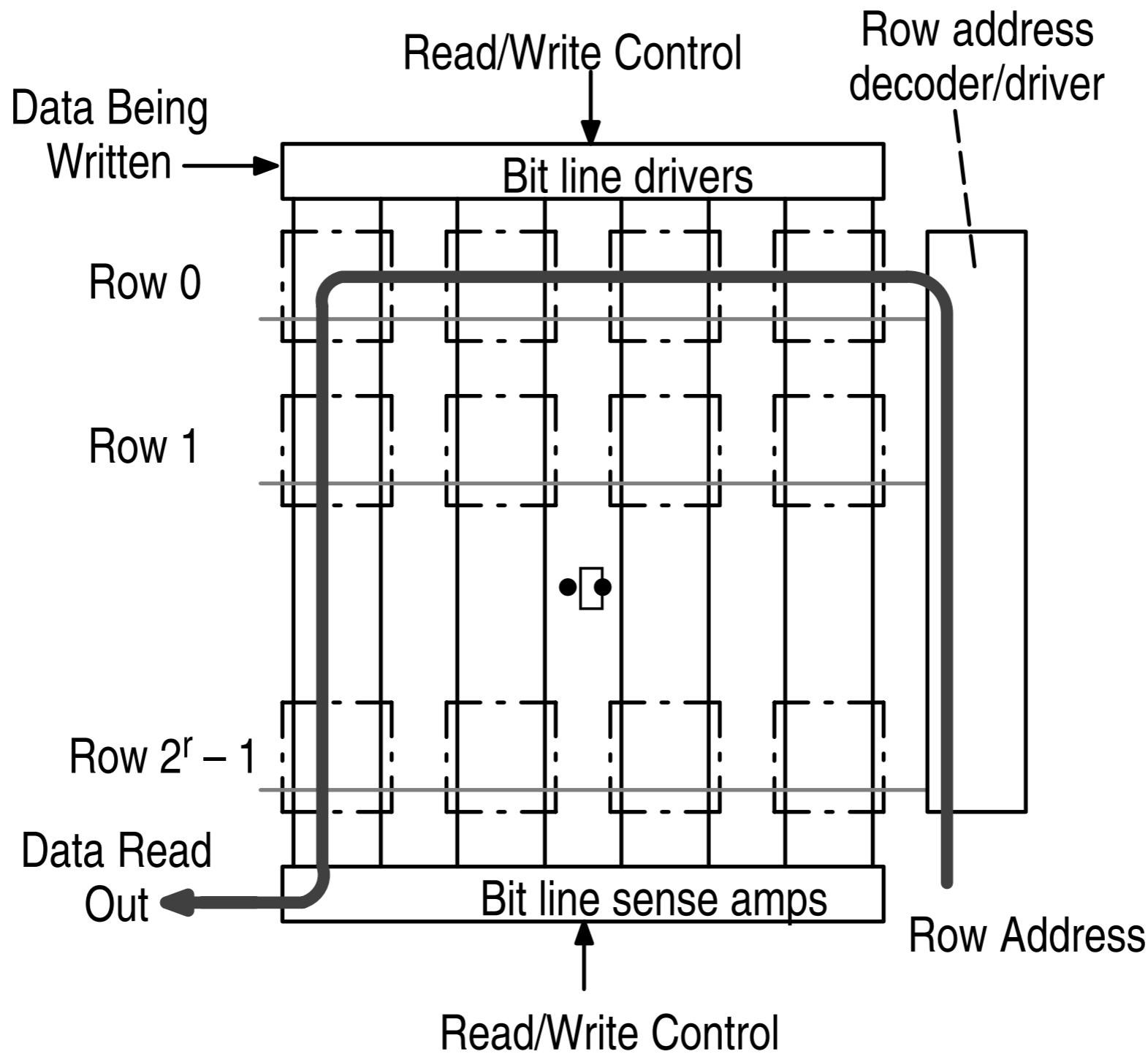


: valid information (1 or 0) on one or more lines

SRAM Internal Organization



- What are the components of the SRAM access time?



- The above figure shows the longest delay path in a SRAM (in gray)

- Between the application of an address and control signals and the time the data is read out, the delays in the path shown are:

$t_{decoder}$ = the delay of the decoder and the driver of the word lines:
this is proportional to the # of address bits (a)

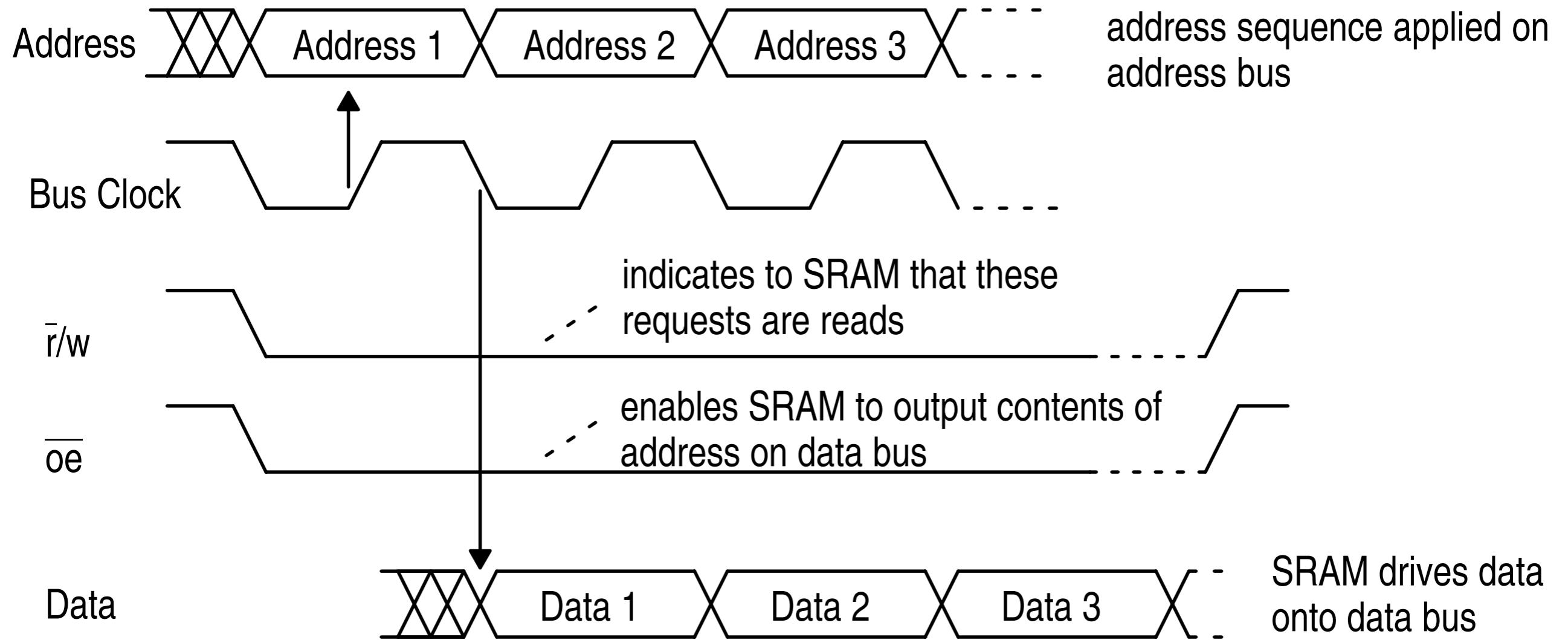
t_{row_line} = the delay of the row wire – from the decoder to the farthest end: this is proportional to the # of bits/word (d)

t_{bit_line} = the delay of the bit lines – from the row farthest from the sense amp to the sense amp: this is proportional to the number of words (2^a)

t_{sense} = the delay of the sense amp

- The read access time is thus ($t_{decoder} + t_{row_line} + t_{bit_line} + t_{sense}$)
 - Note how the access time depends on the dimensions of the array and why “smaller is faster” holds in this case.
 - Possible ways to speed up the SRAM access time:
 - (a) Reduce the array dimensions
 - (b) Use smaller transistors (i.e., a more aggressive fabrication process): this reduces the various proportionality constants.

Pipelining SRAM access



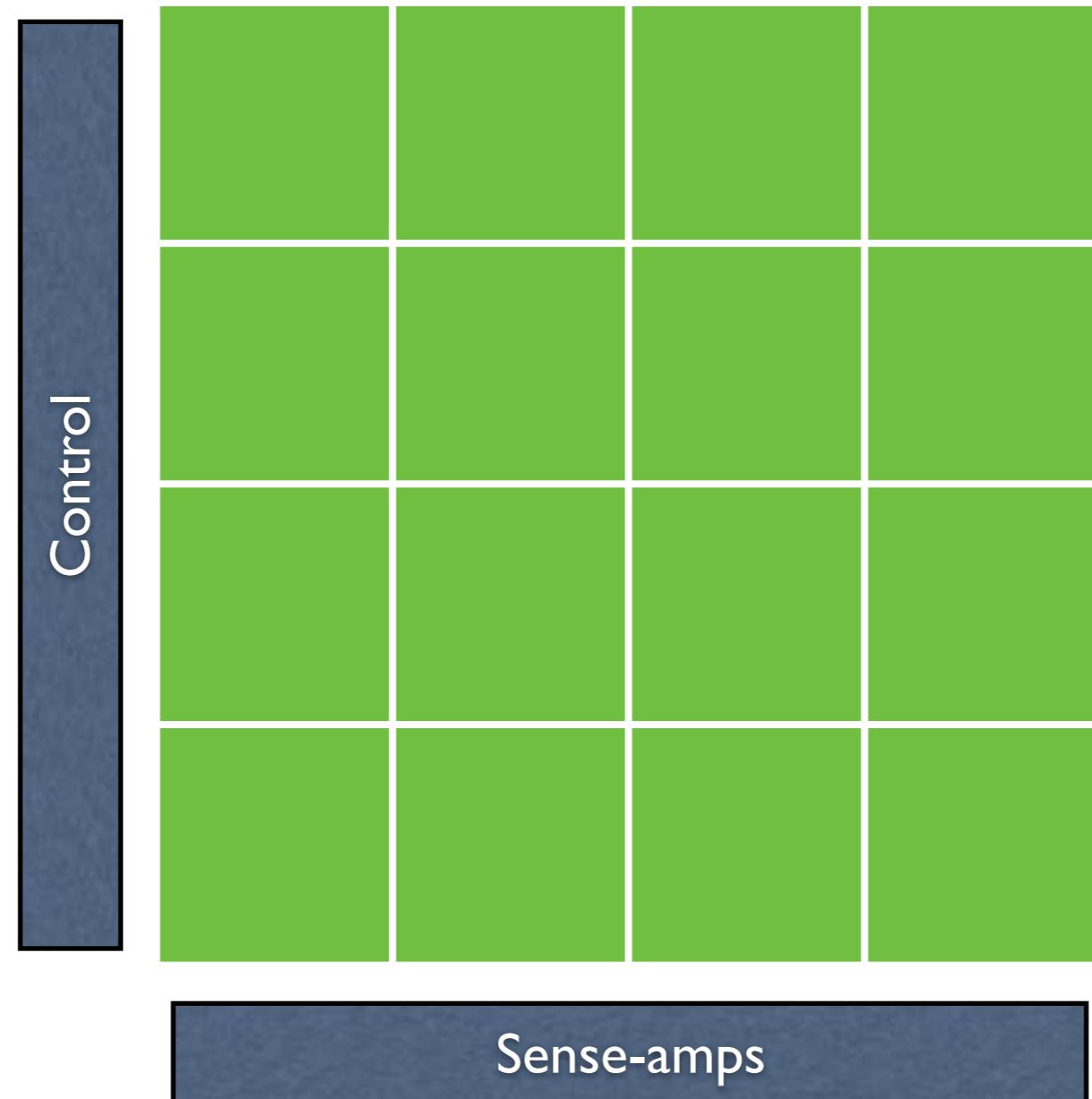
What do we make from SRAMs?

- Physical register files
- Register alias tables
- Reorder buffers
- Scoreboards
- Instruction Queues

Register files and ILP

- Superscalar processors perform many simultaneous accesses to a register file
- More ports requires more area
- SRAM delay is a function of area

Sub-banking (vertical stripes)



Multi-banking (horizontal stripes)



Sense-amps

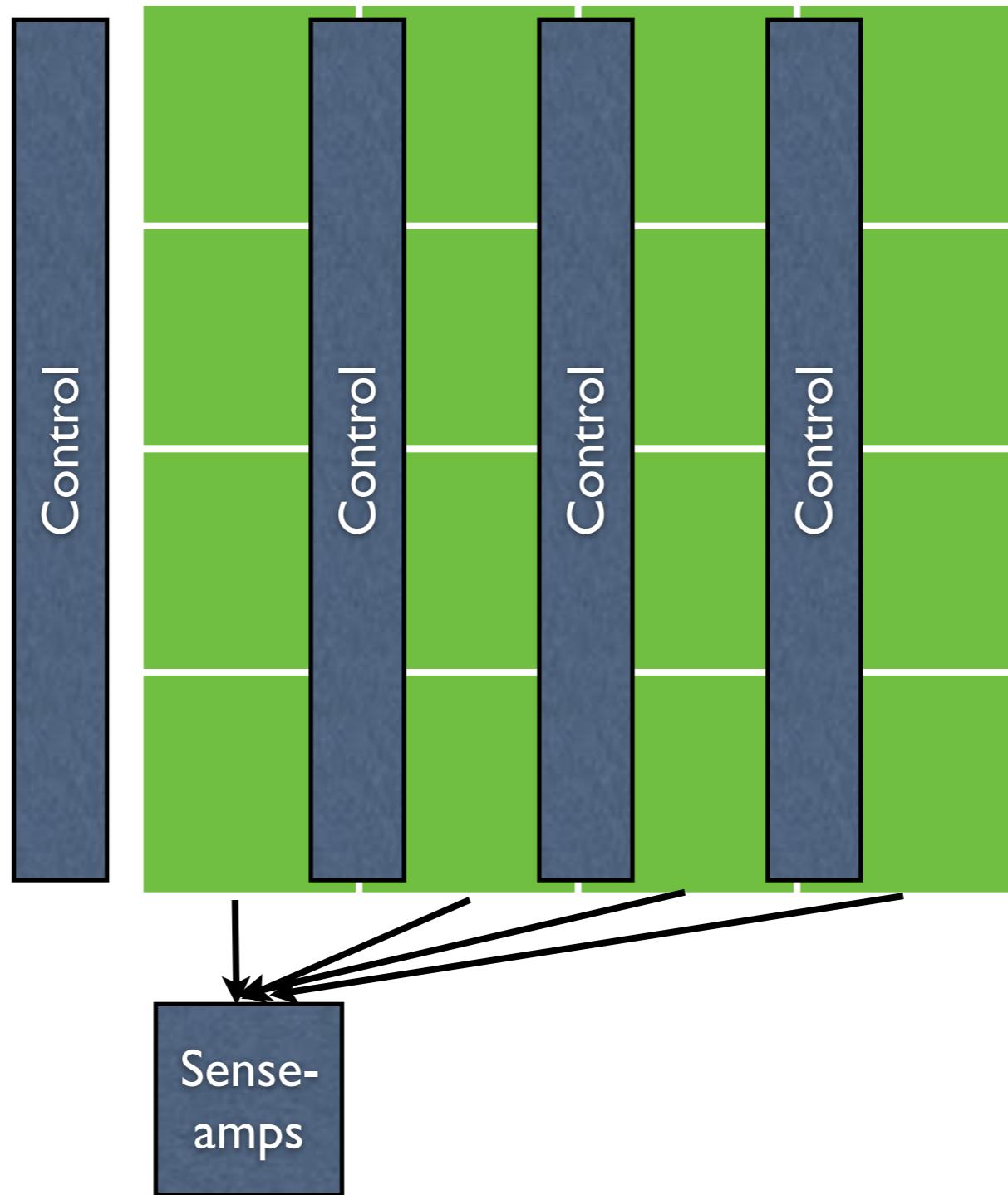


Sense-amps



Sense-amps

Sub-banking (vertical stripes)



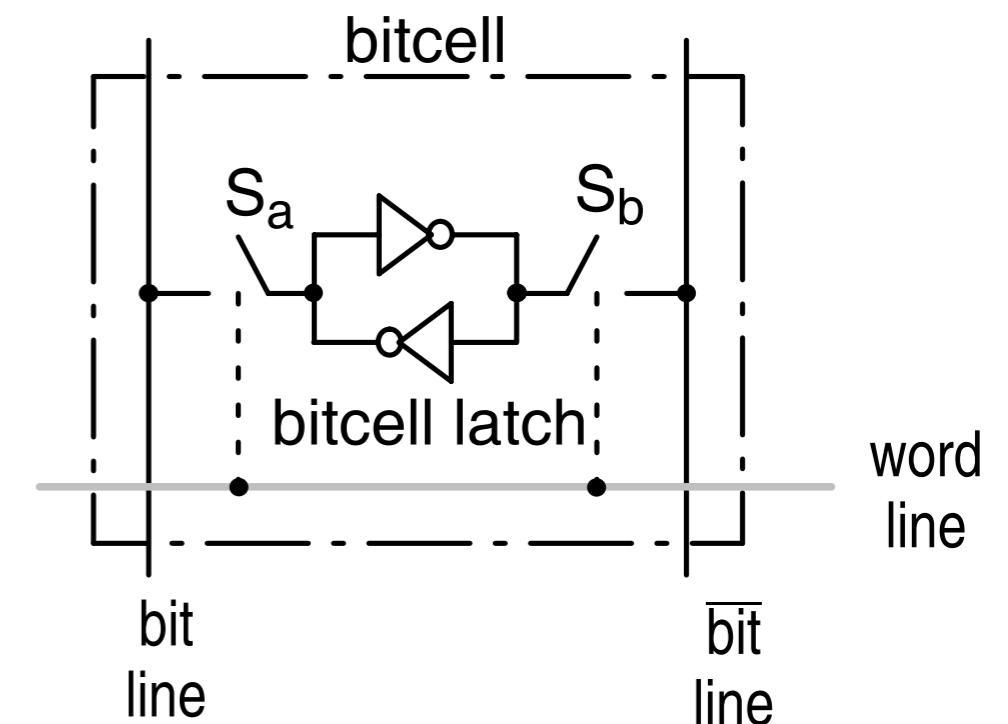
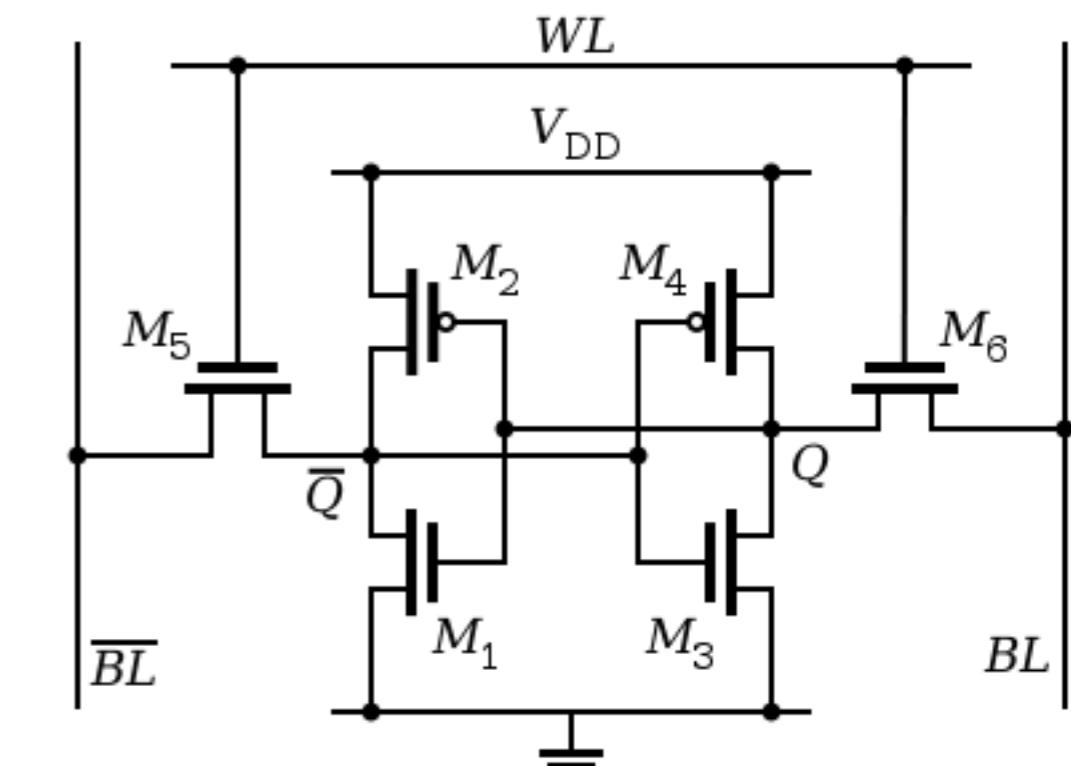
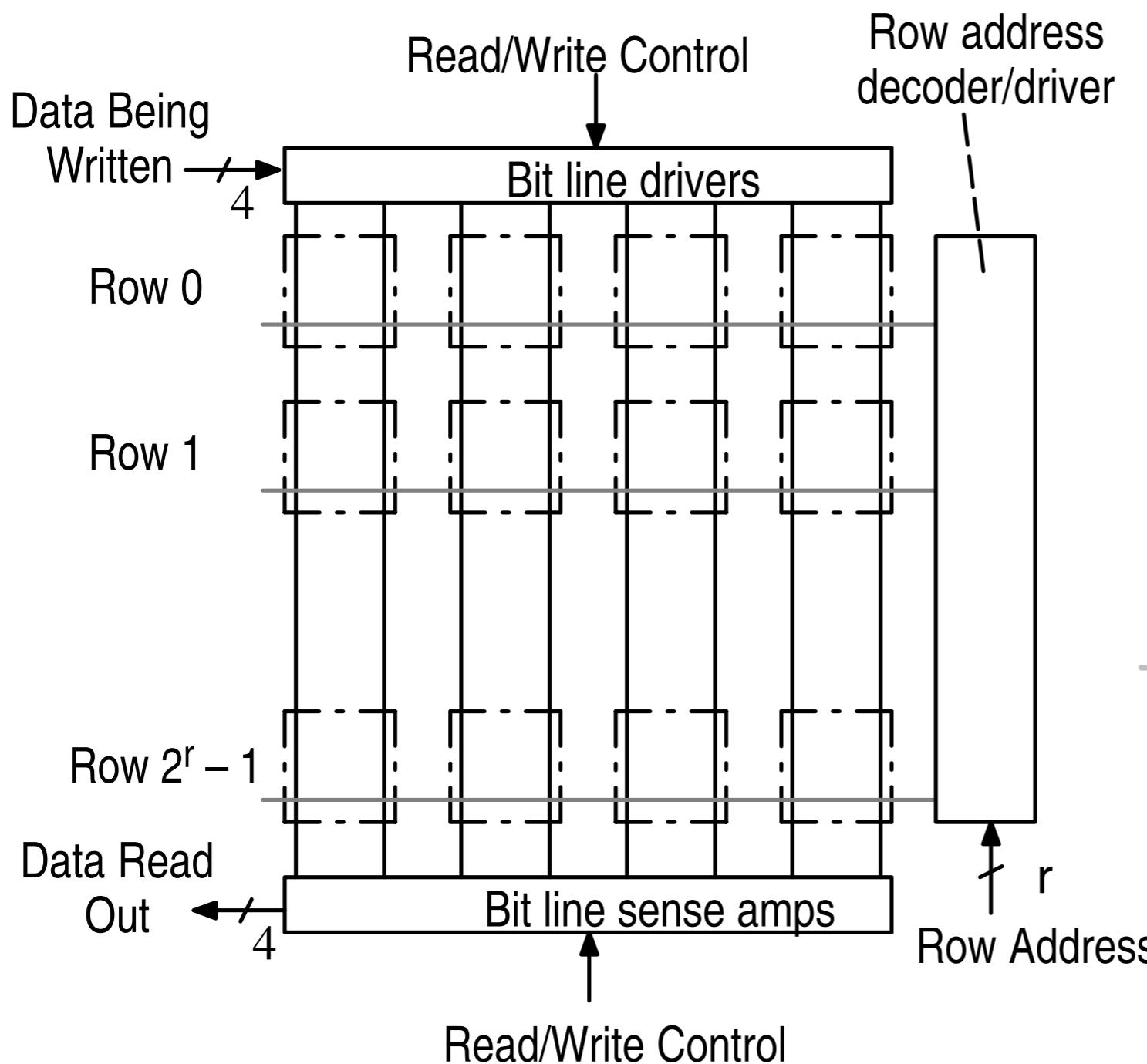
Clustered



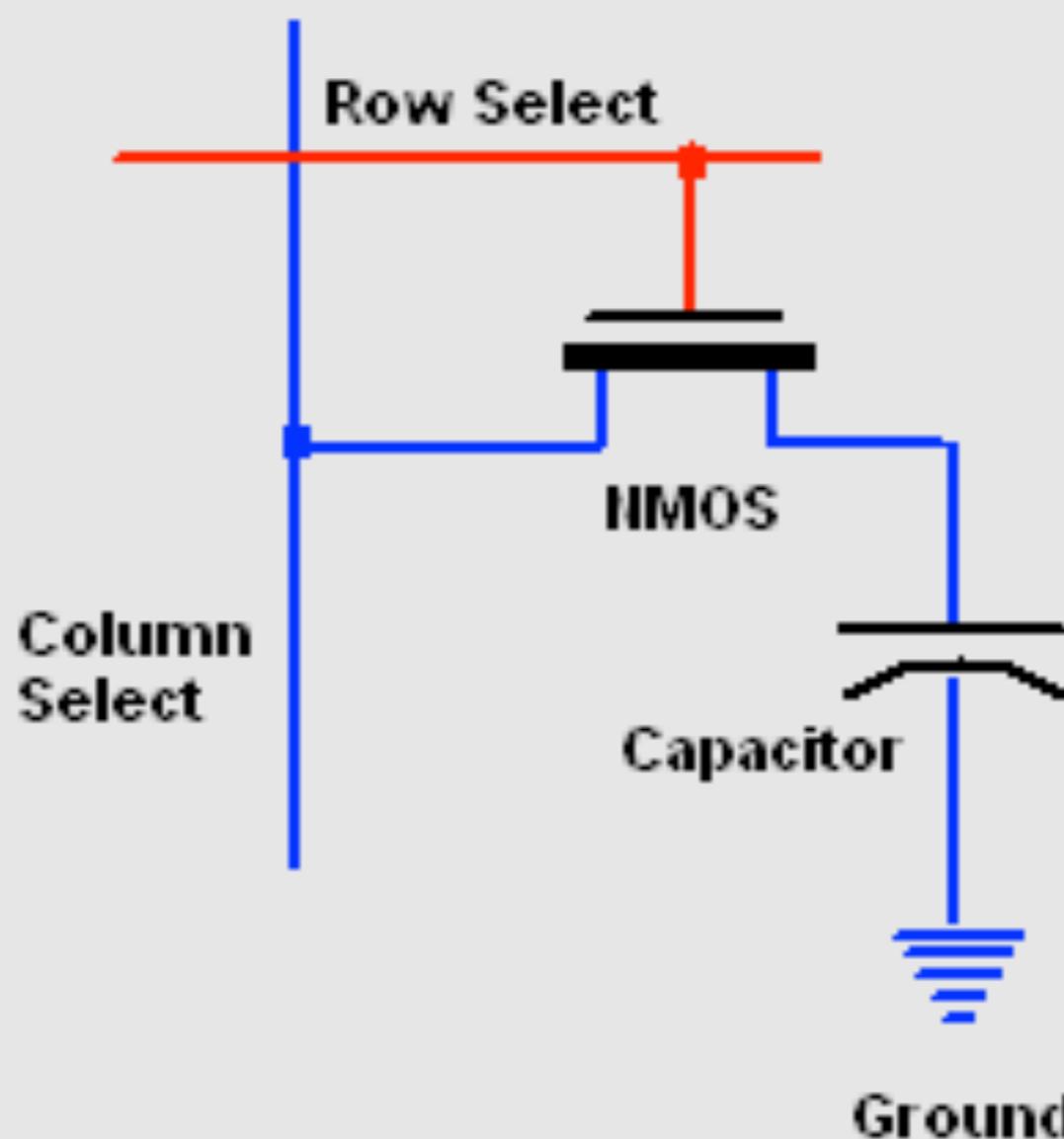
DRAM Organization

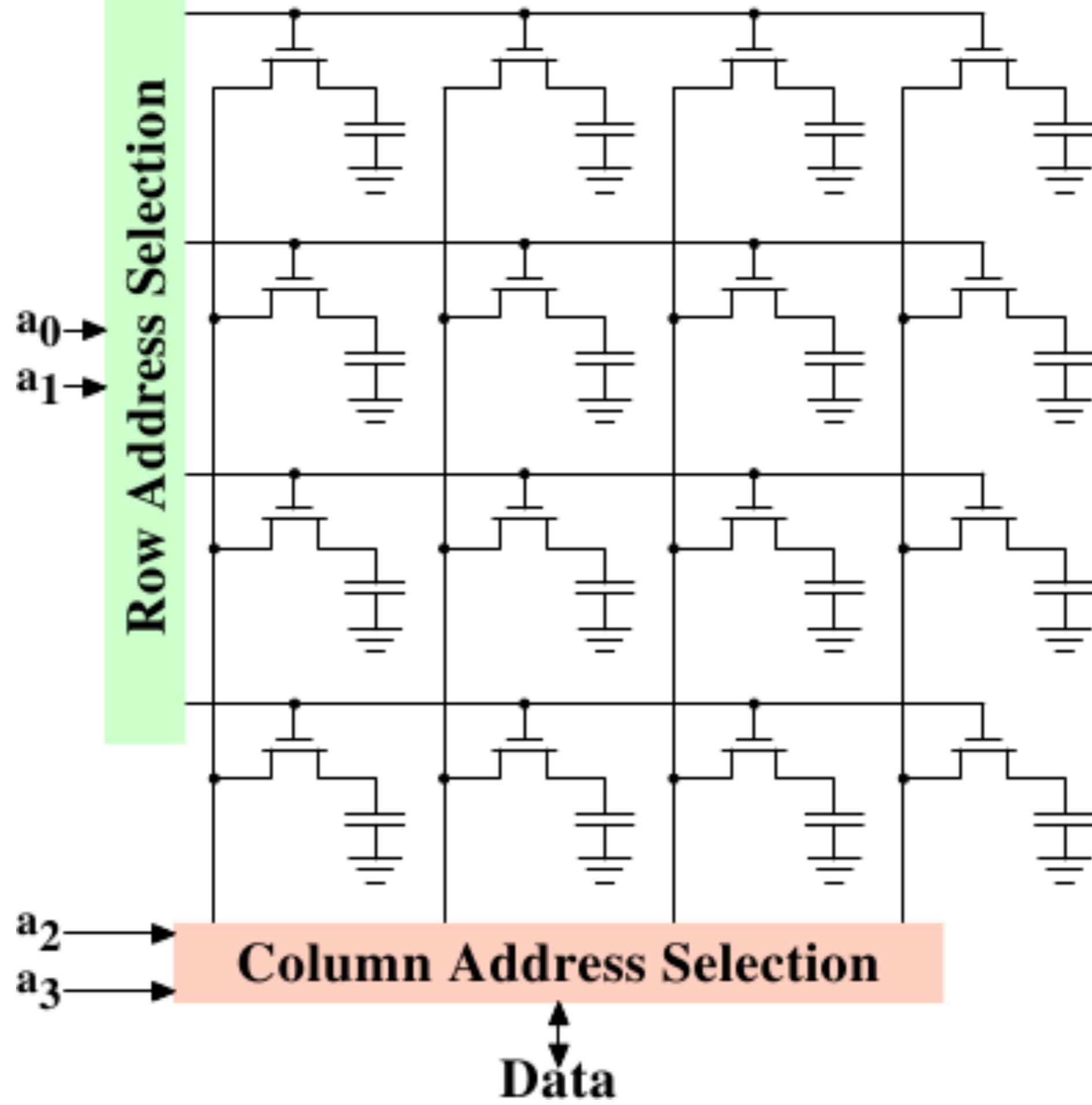
- Stored as 2D array with very long rows (e.g. 8K bits)
- “Activating” copies row into row buffer
- Reads and write columns in row buffer
- “Precharge” copies row back into array

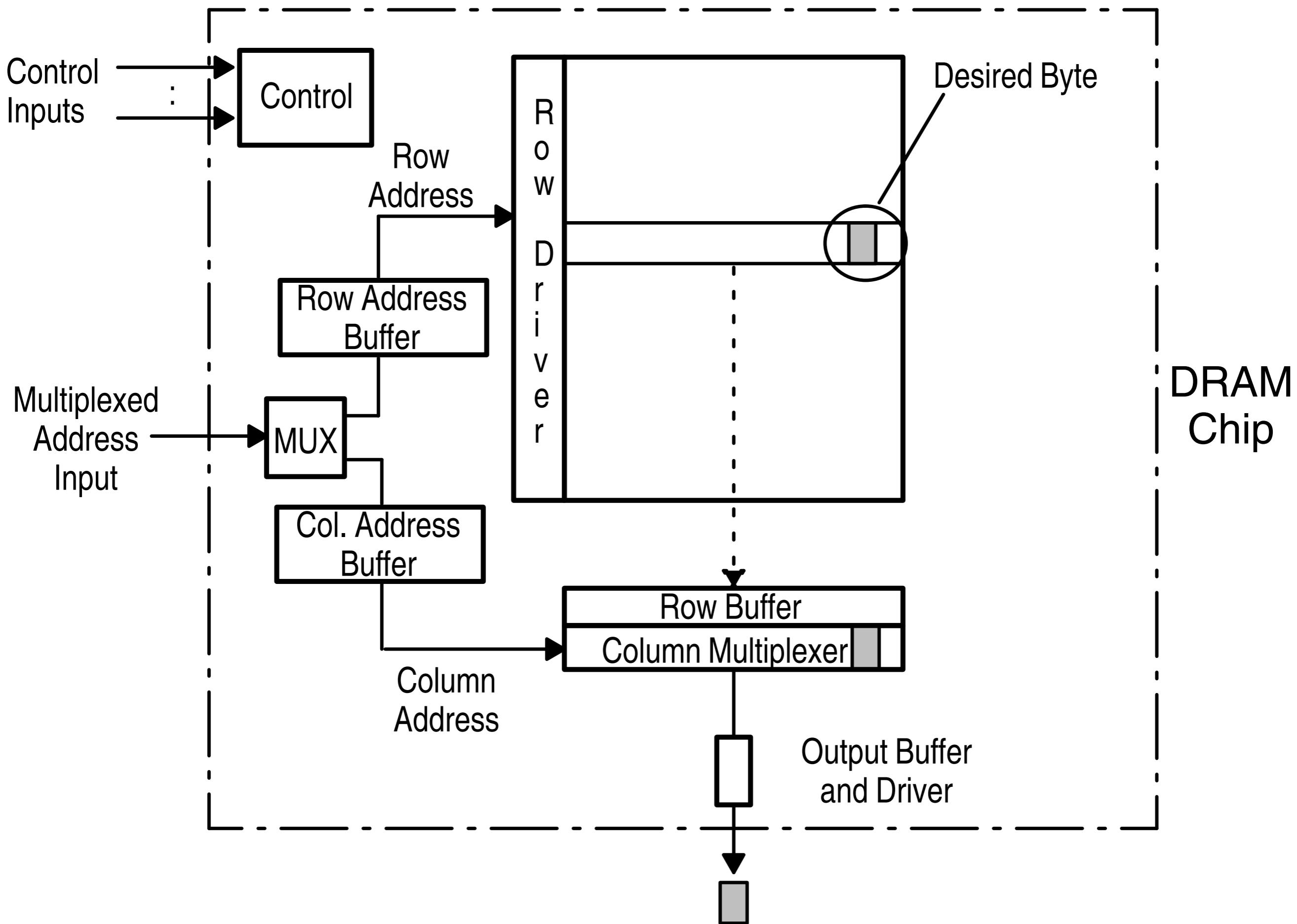
SRAM Internal Organization

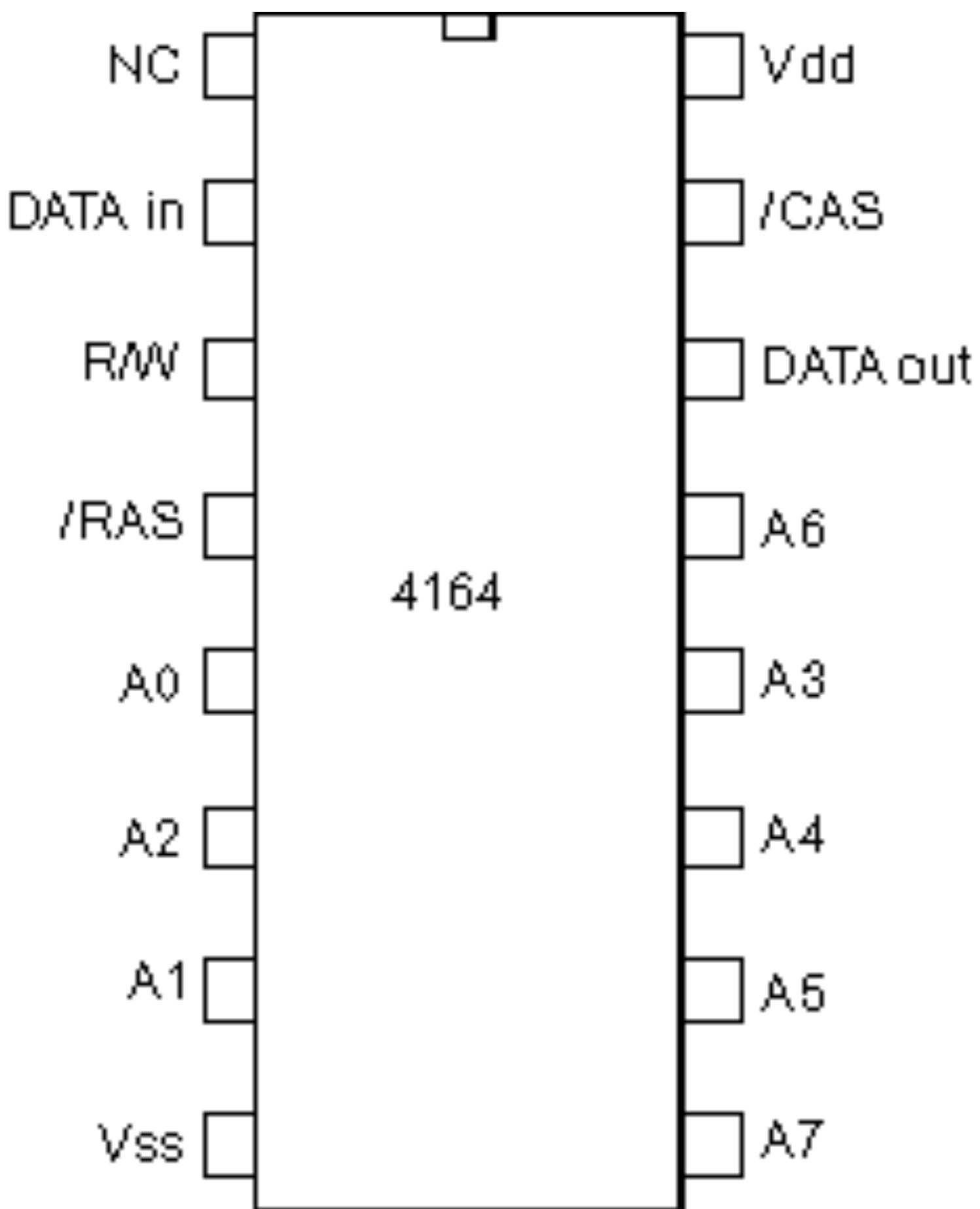


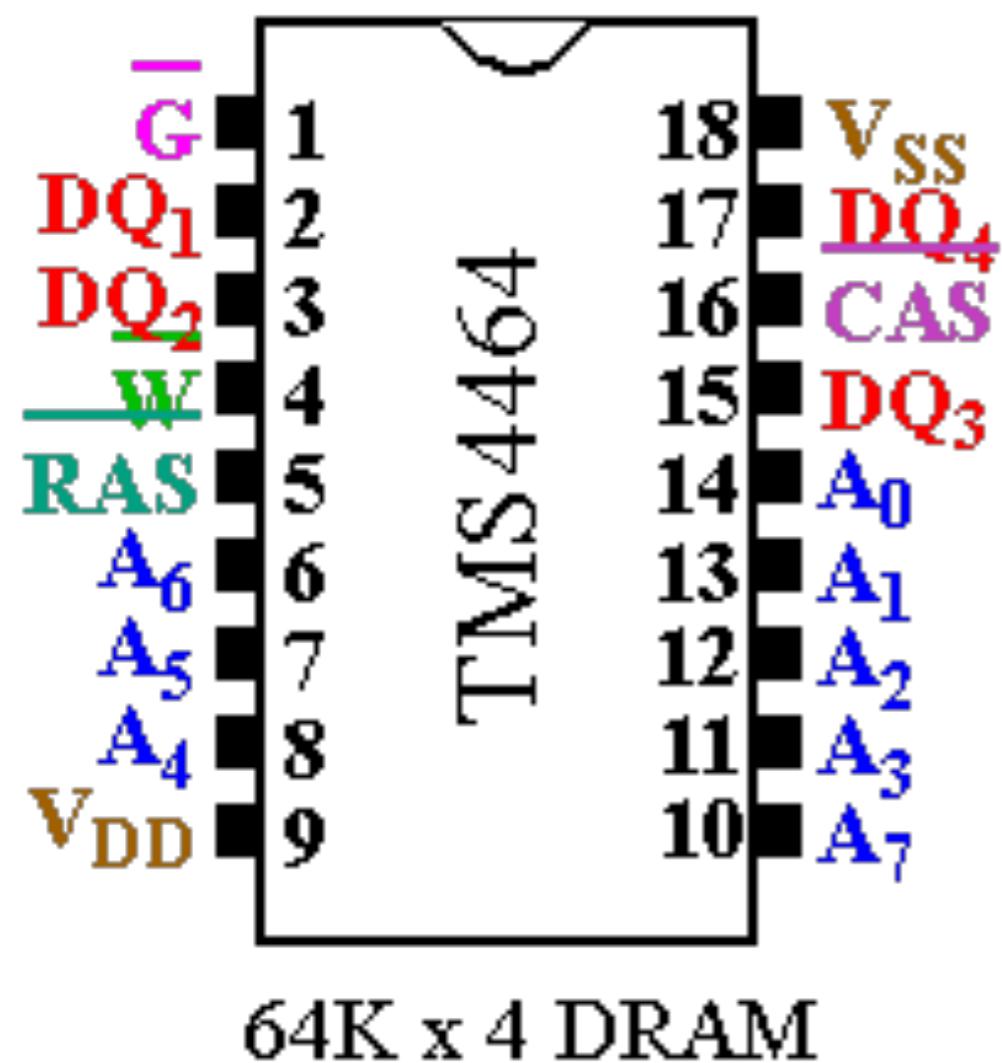
One Dynamic RAM Bit (DRAM bit)



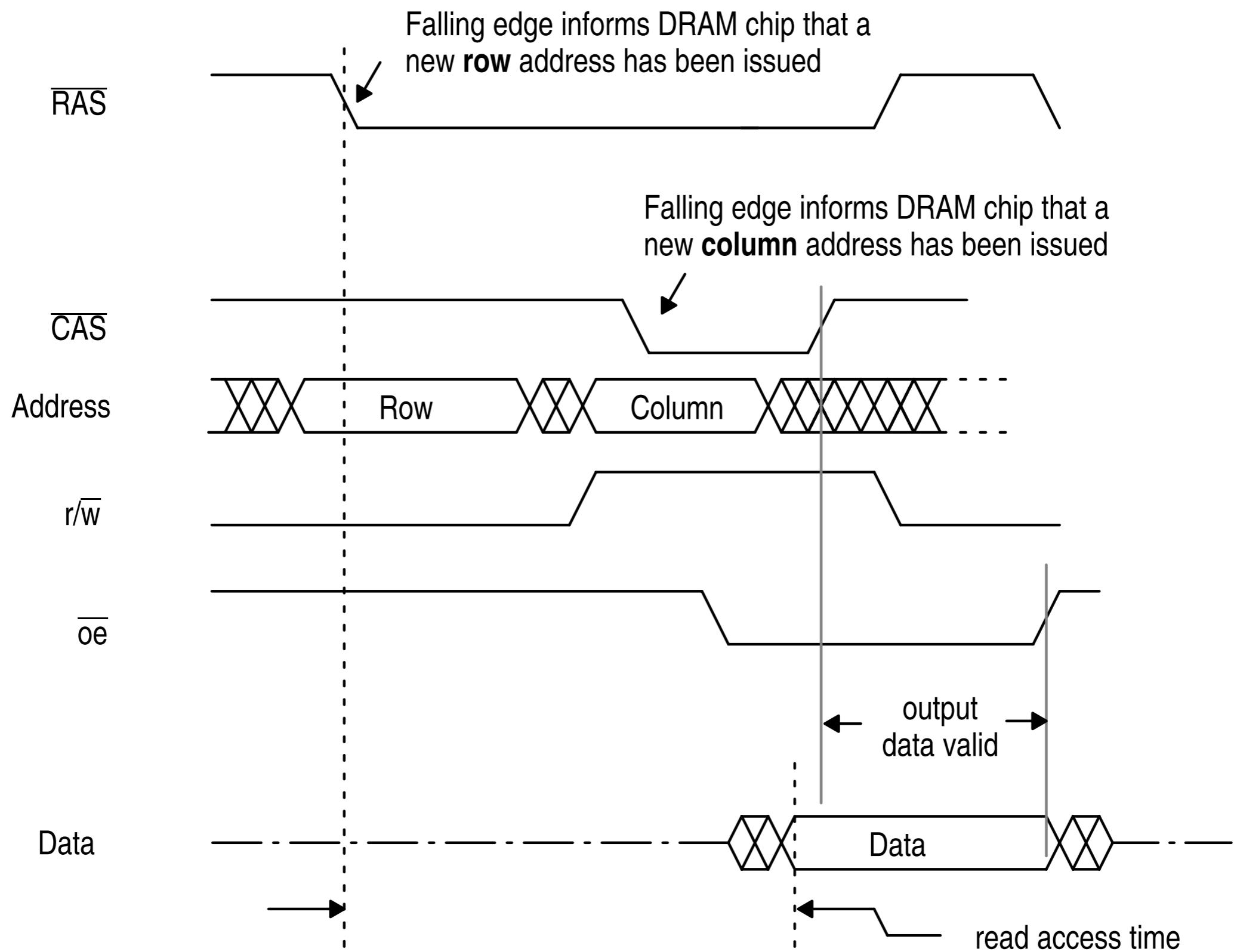








Pin(s)	Function
A ₀ -A ₇	Address
DQ ₀ -DQ ₄	Data In/Data Out
RAS	Row Address Strobe
CAS	Column Address Strobe
G	Output Enable
W	Write Enable



= undefined



: valid information (1 or 0) on one or more lines

— - - - - = output floating (open)

DRAM Facts

- DRAM Bit cells **leak**
 - Require refresh every 10 to 15 microseconds
- A DRAM controller:
 - Must perform complex scheduling
 - Generates read, write, activate, precharge, and refresh commands
 - Schedules refresh
 - Performs error checking

Refresh Mechanism

- Read contents of row into row buffer before contents degrade
- Write row back into DRAM array
 - Replenishes charge in cells holding ‘1’

Refresh Cycles

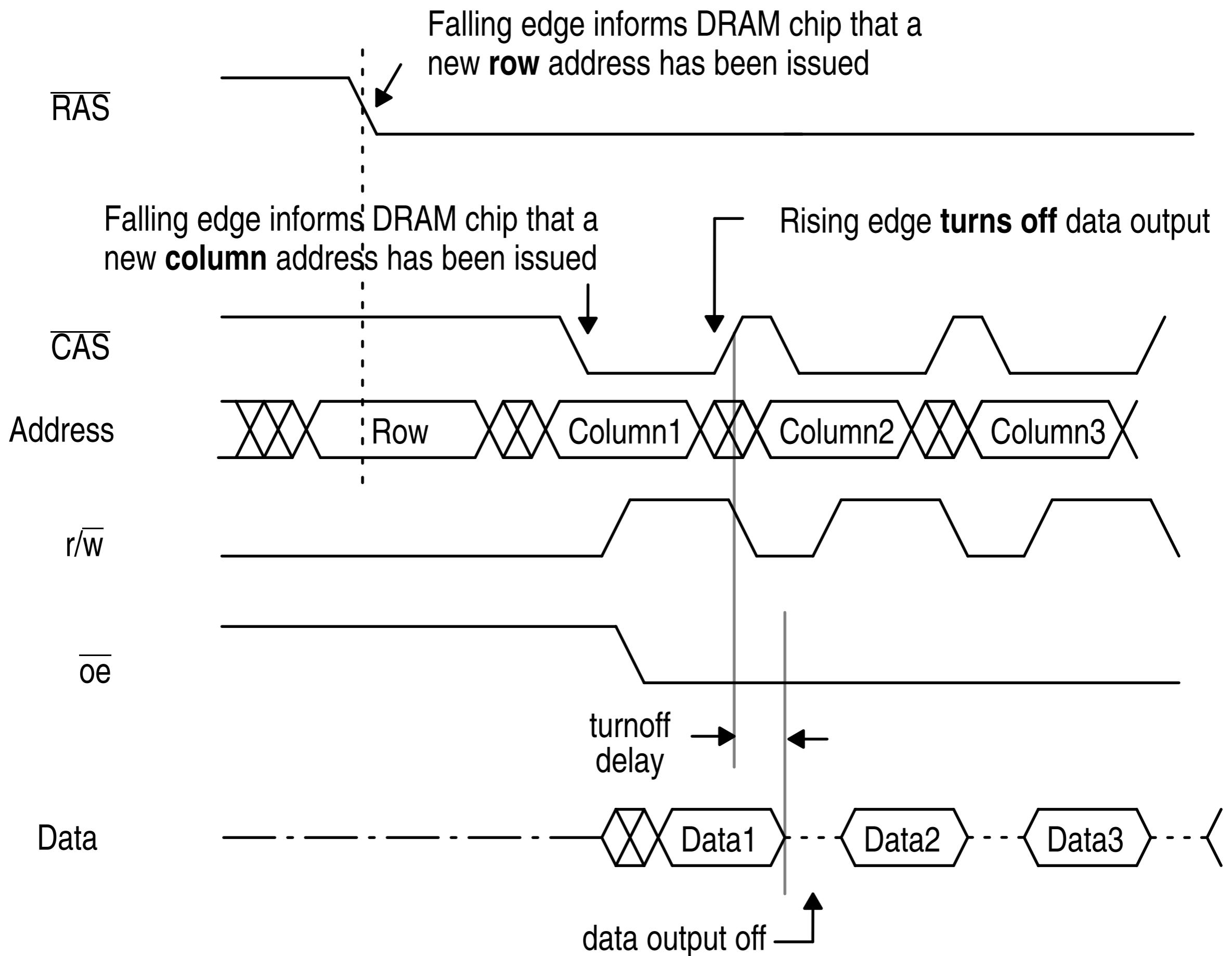
- Distributed mode -- Evenly spreads out refreshed amid other activity
- Burst mode -- refresh many rows consecutively during idle periods

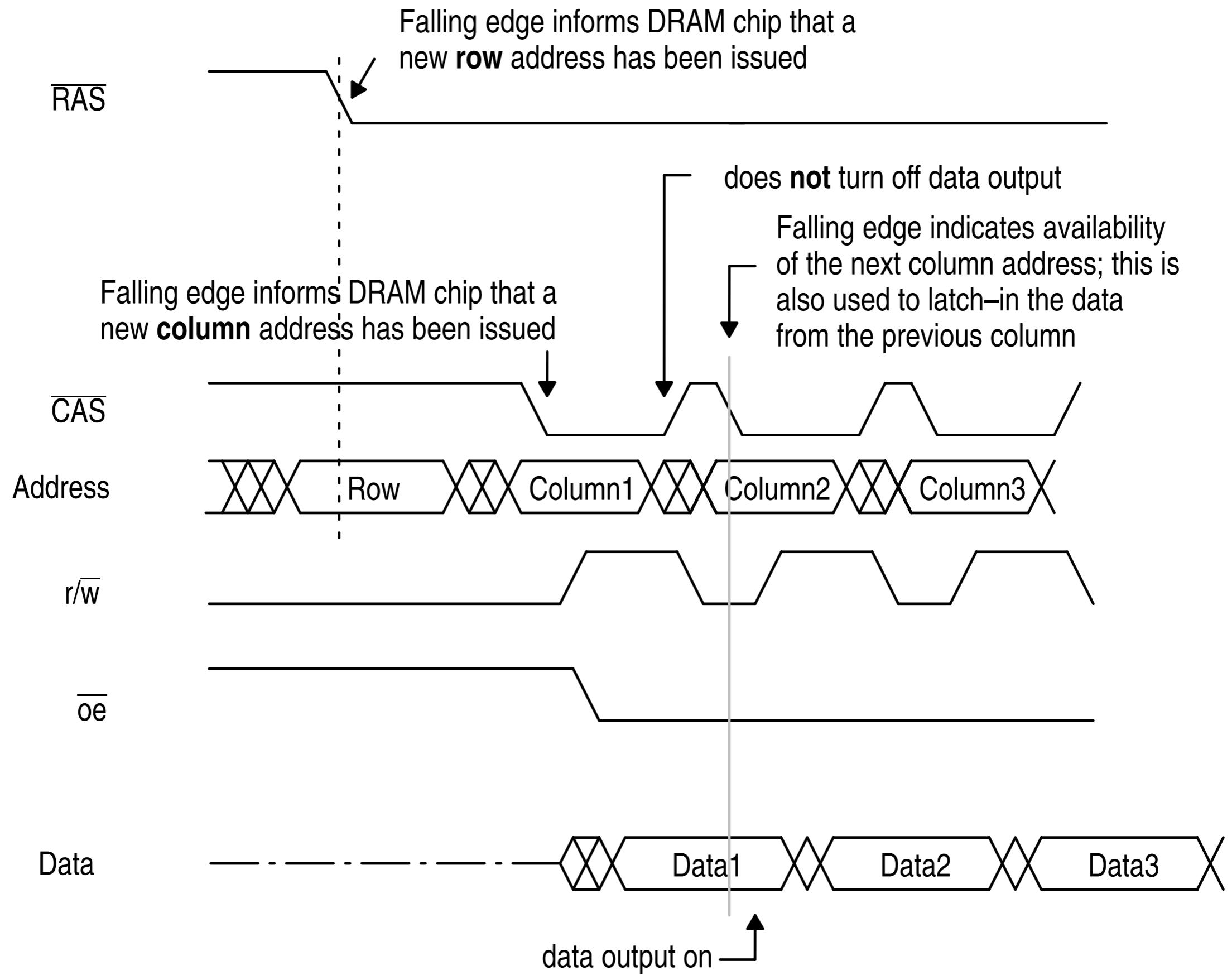
Refresh Options

- RAS-only
- CAS-before-RAS
- Hidden

DRAM Interface Evolution

- Asynchronous interfaces
- Fast page mode (FPM) -- allows multiple column access to same row
- Extended Data Out (EDO) -- primitive pipelining
- Synchronous DRAM -- fully pipelined





Synchronous DRAM

- Interface is now synchronous, timed by an explicit clock signal
- Interface is pipelined, allowing for high clock rates and high throughput for column access
- Introduces “CAS latency”
- With DDR, data bus faster than control bus
- More recent advancements focus on signal integrity

Multiple Banks

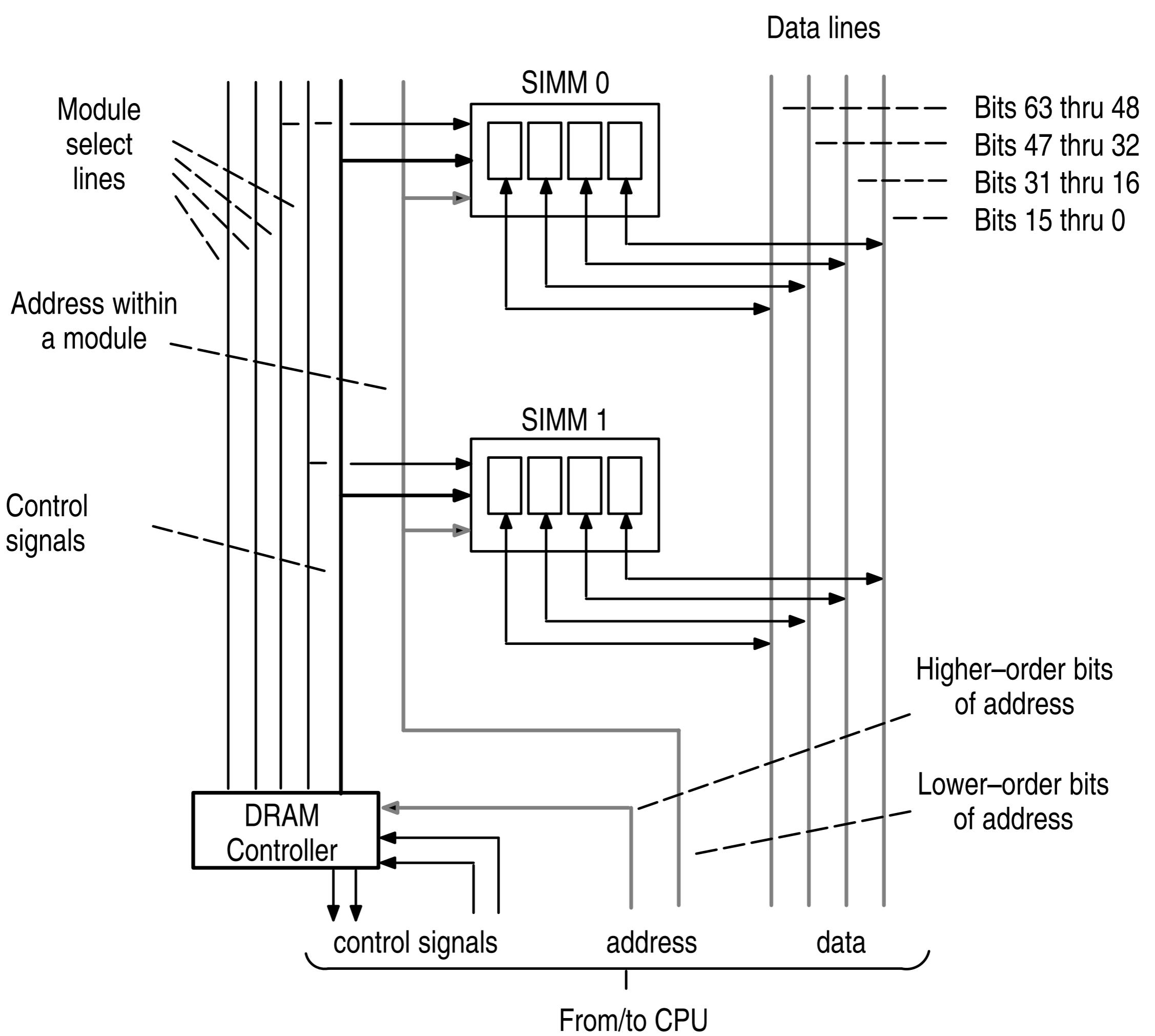
- DRAM internals divided into multiple independent banks
- Each has its own row buffer
- Share common interface
- Used to hide activate/precharge latencies

DRAM Commands

- RAS-CAS-WE
- 111 - No operation
- 110 - Burst terminate
- 101 - Read (CAS low, WE high)
- 100 - Write (CAS low, WE low)
- 011 - Activate (RAS low)
- 010 - Precharge
- 001 - Auto-refresh
- 000 - Load mode register

DRAM modules

- DRAM chips packages into modules
 - SIMMs and DIMMs
- 8 chips on one module constitute a “rank”
 - Common CS, all access in parallel
 - Controller has one CS for each rank

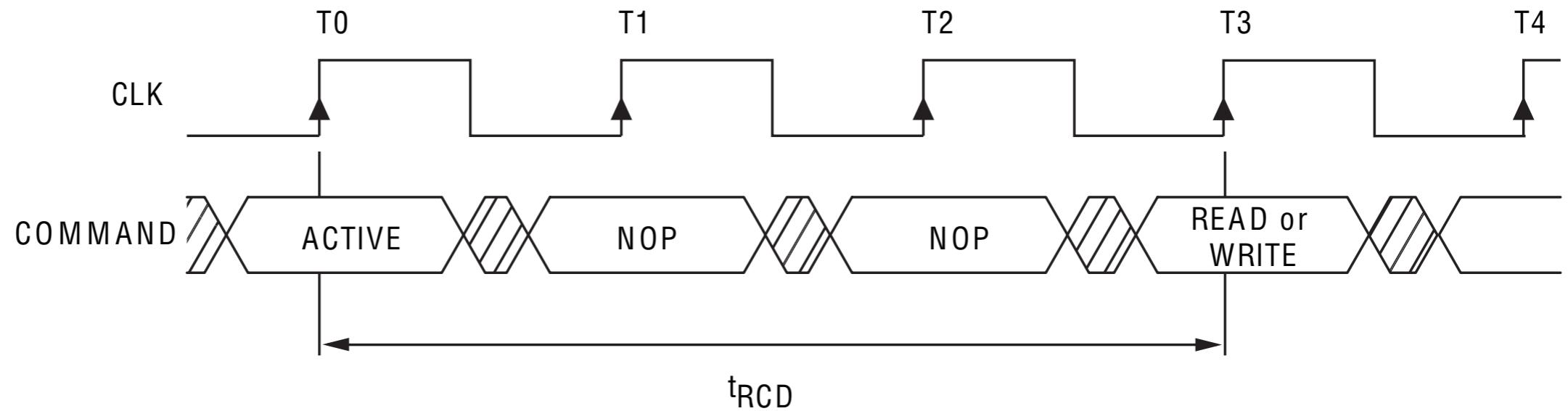


DRAM Timing Numbers

From...	To...	Activate	Read	Write	Precharge
Activate	Activate	tRRD	tRCD	tRCD	tRAS
Read	Read	--	1	tRTW	tRTP
Write	Write	--	tWTR	1	tWR
Precharge	Precharge	tRP	--	--	tPTP

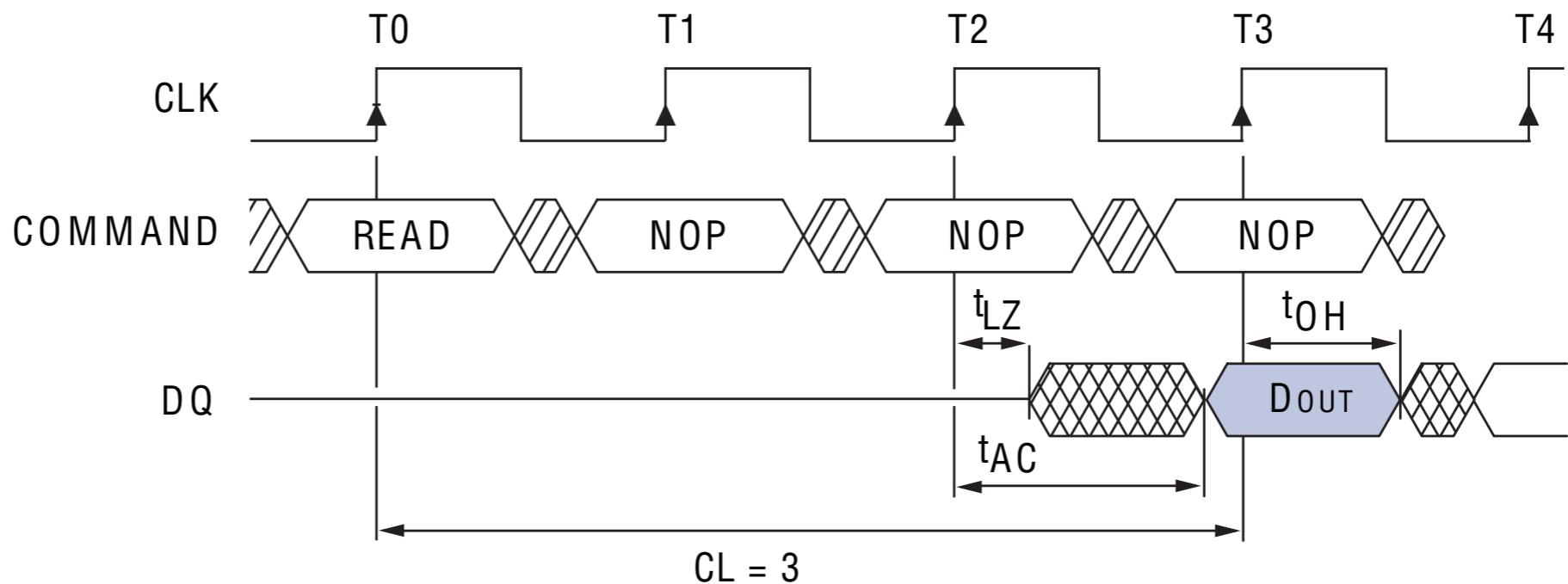
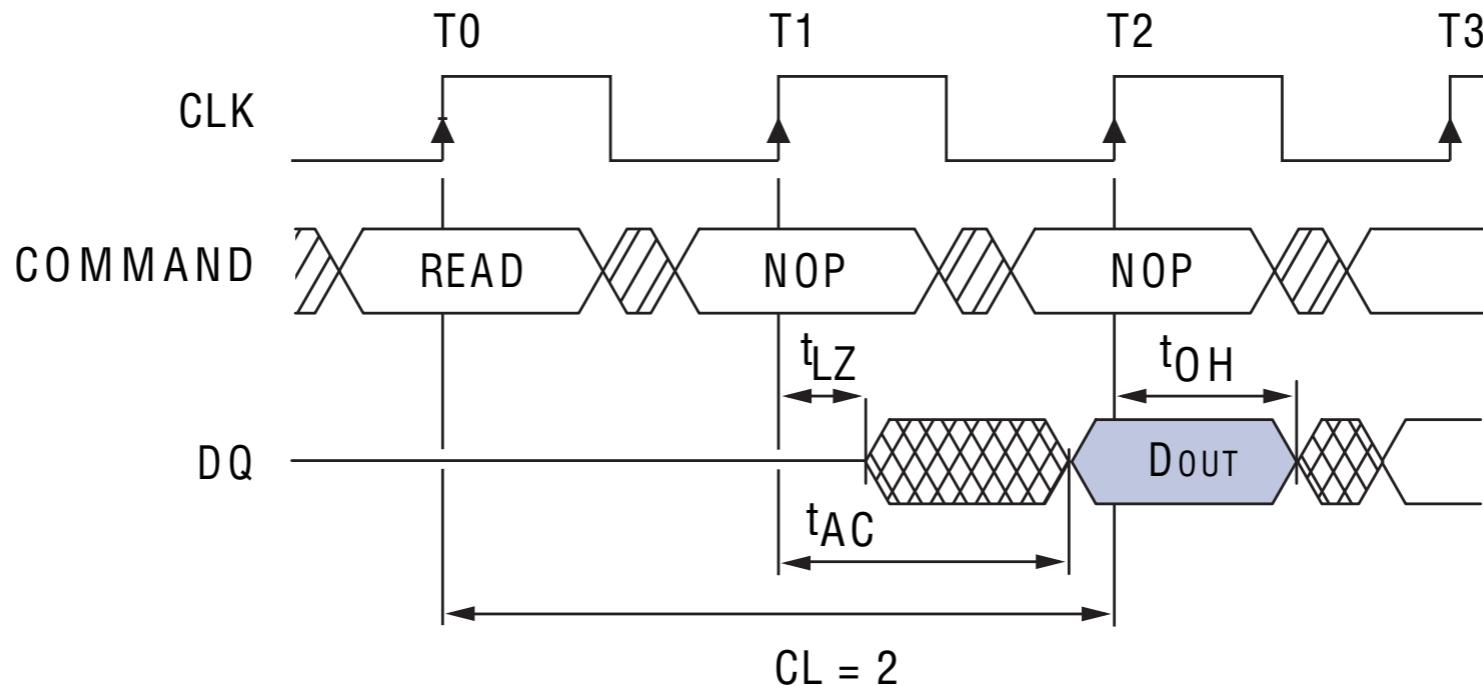
tCL = pipeline latency for reads

Example Meeting $t_{RCD} \text{ (MIN)}$ when $2 < t_{RCD} \text{ (MIN)}/t_{CK} \leq 3$



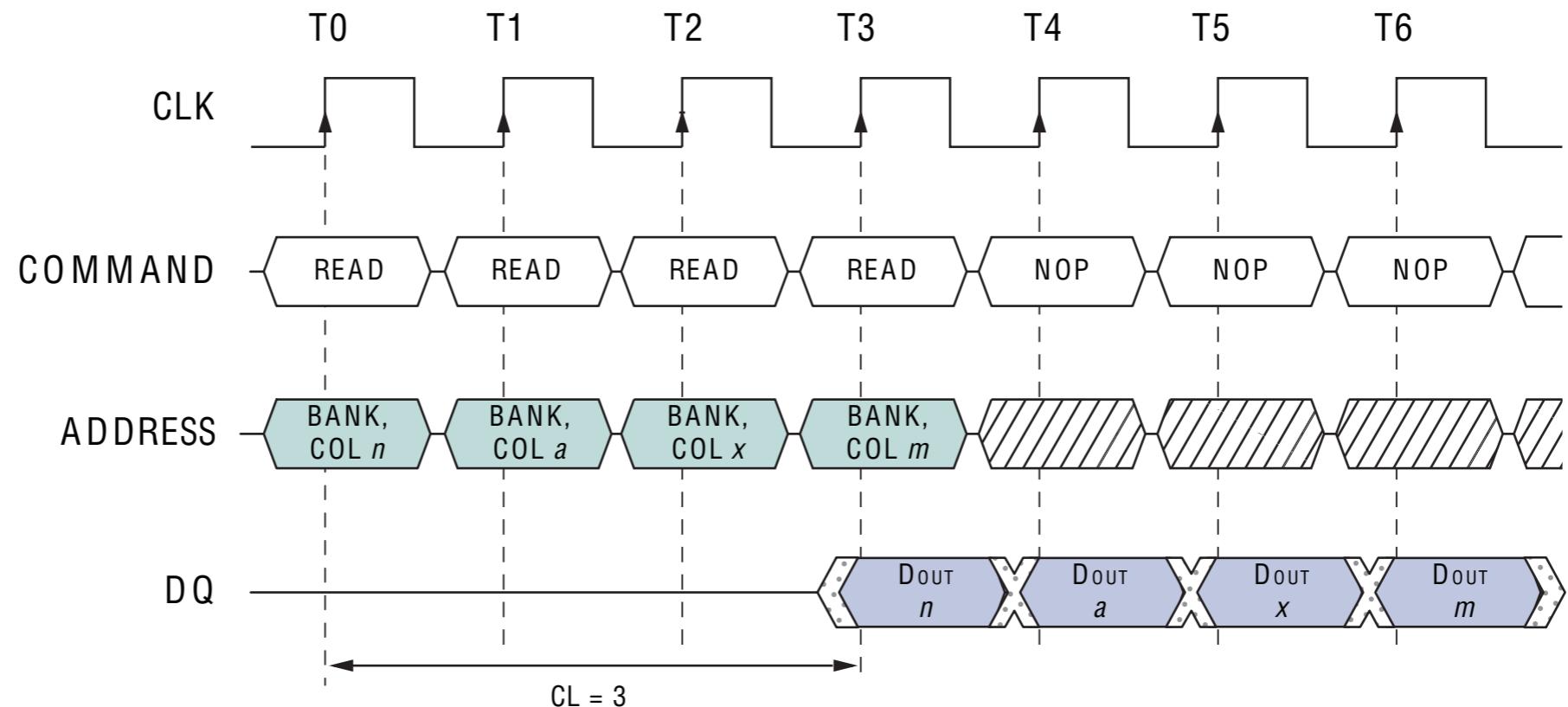
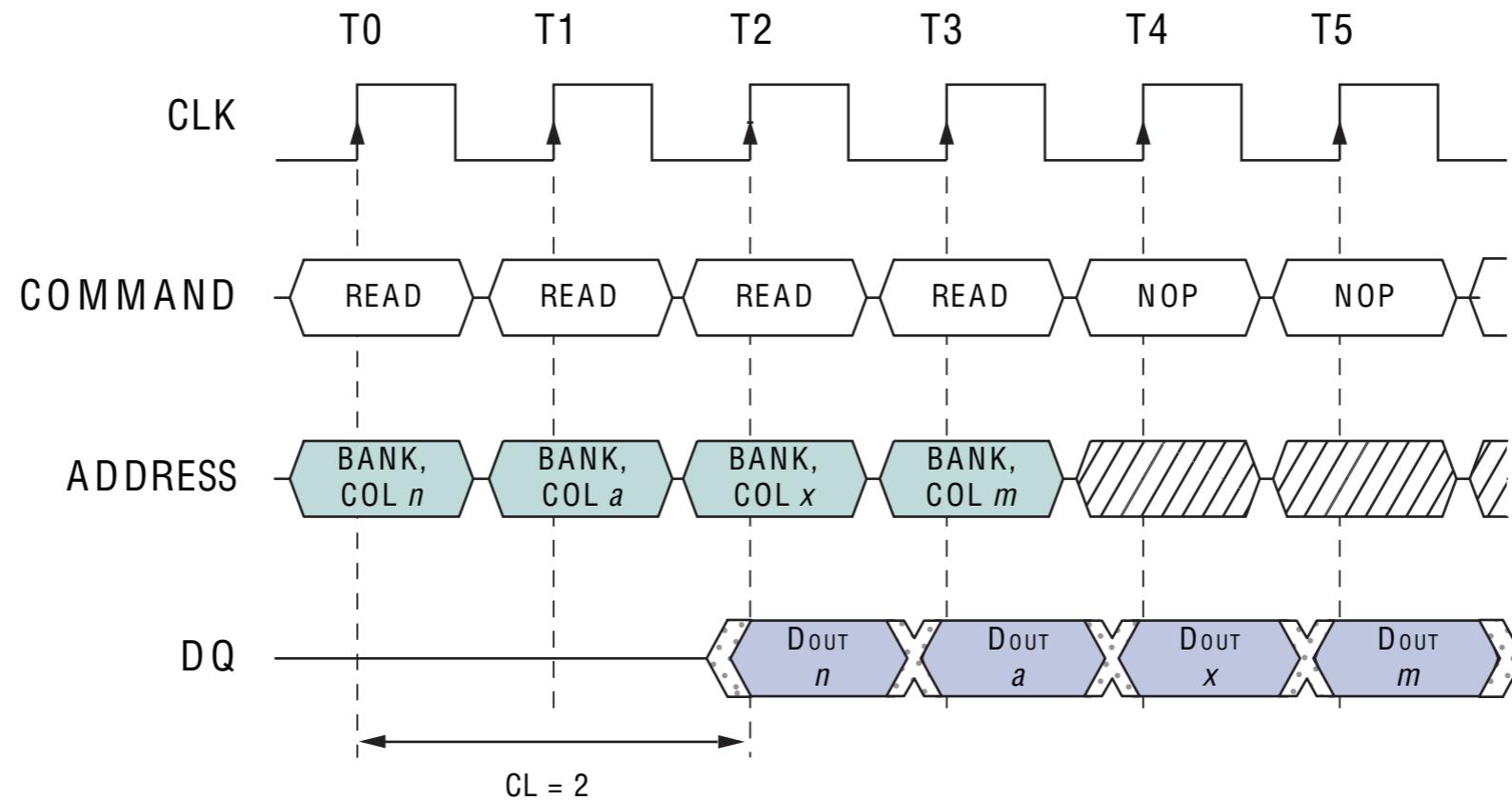
 Don't Care

CAS Latency



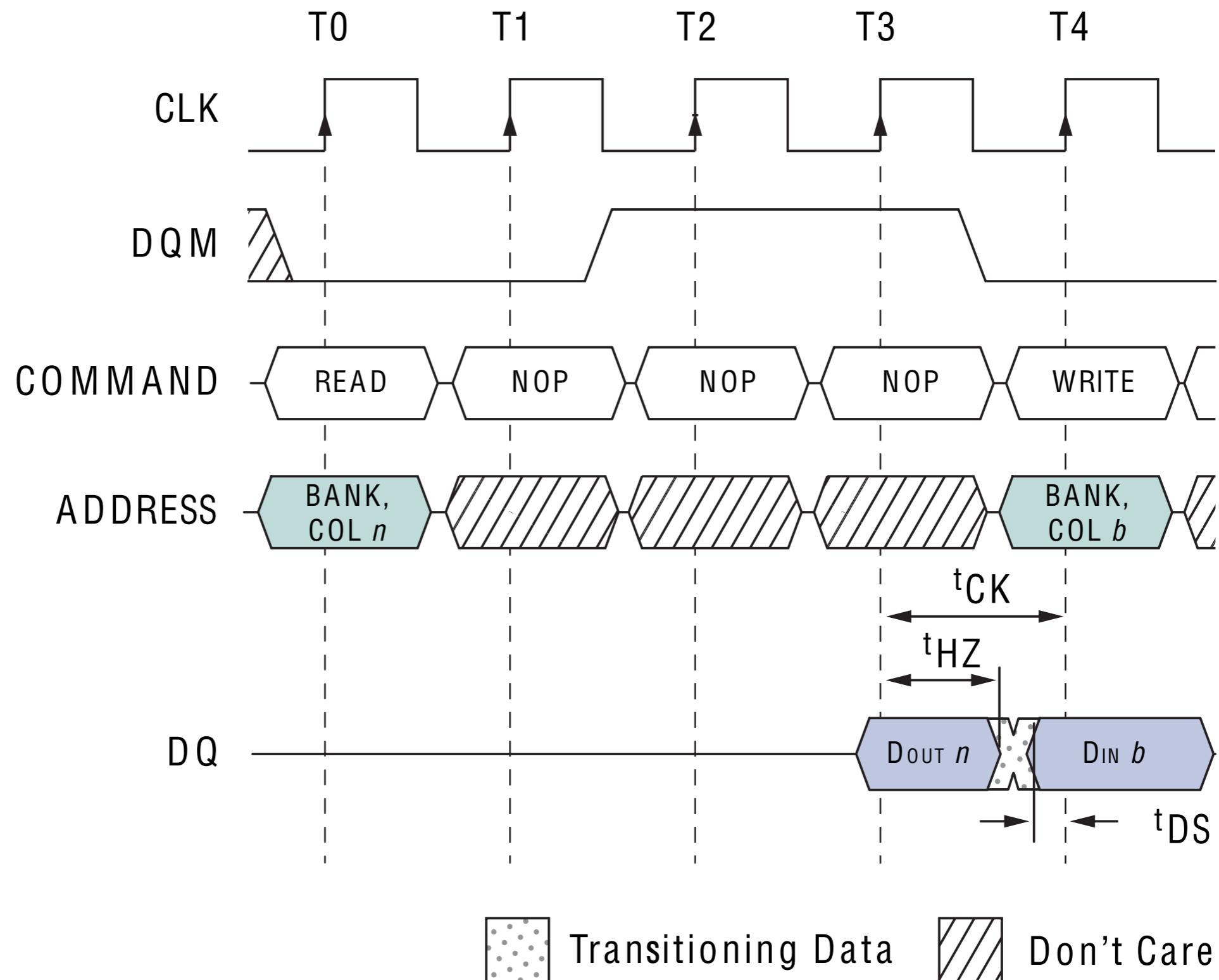
 Don't Care
 Undefined

Random READ Accesses

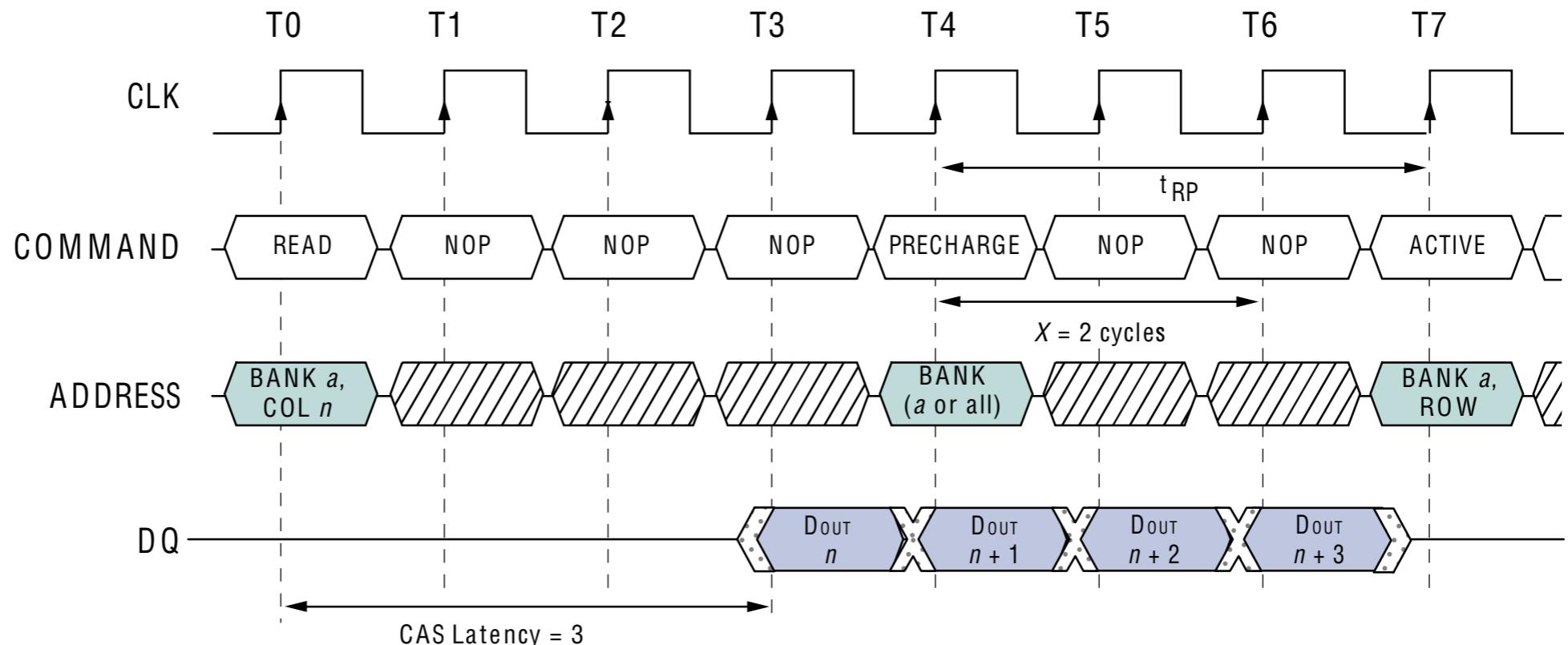
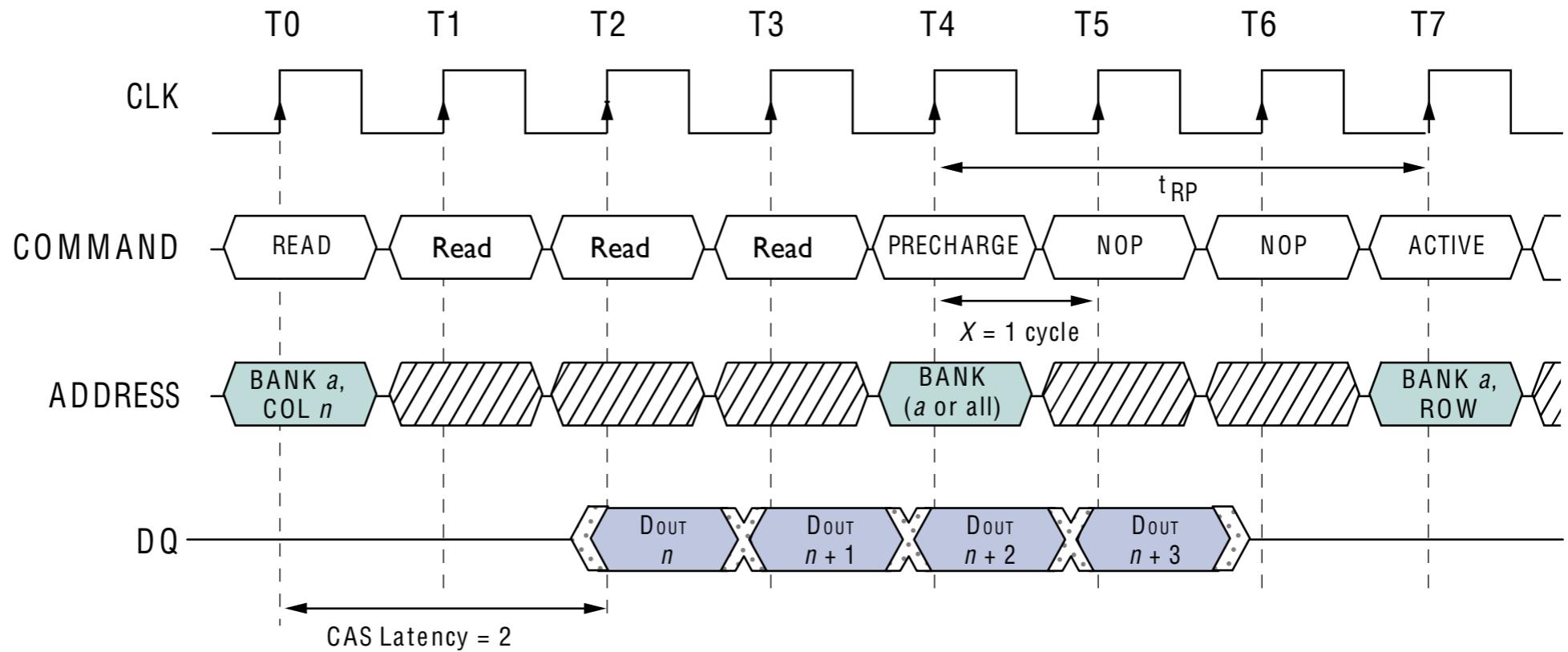


Transitioning Data Don't Care

READ-to-WRITE



READ-to-PRECHARGE

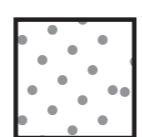
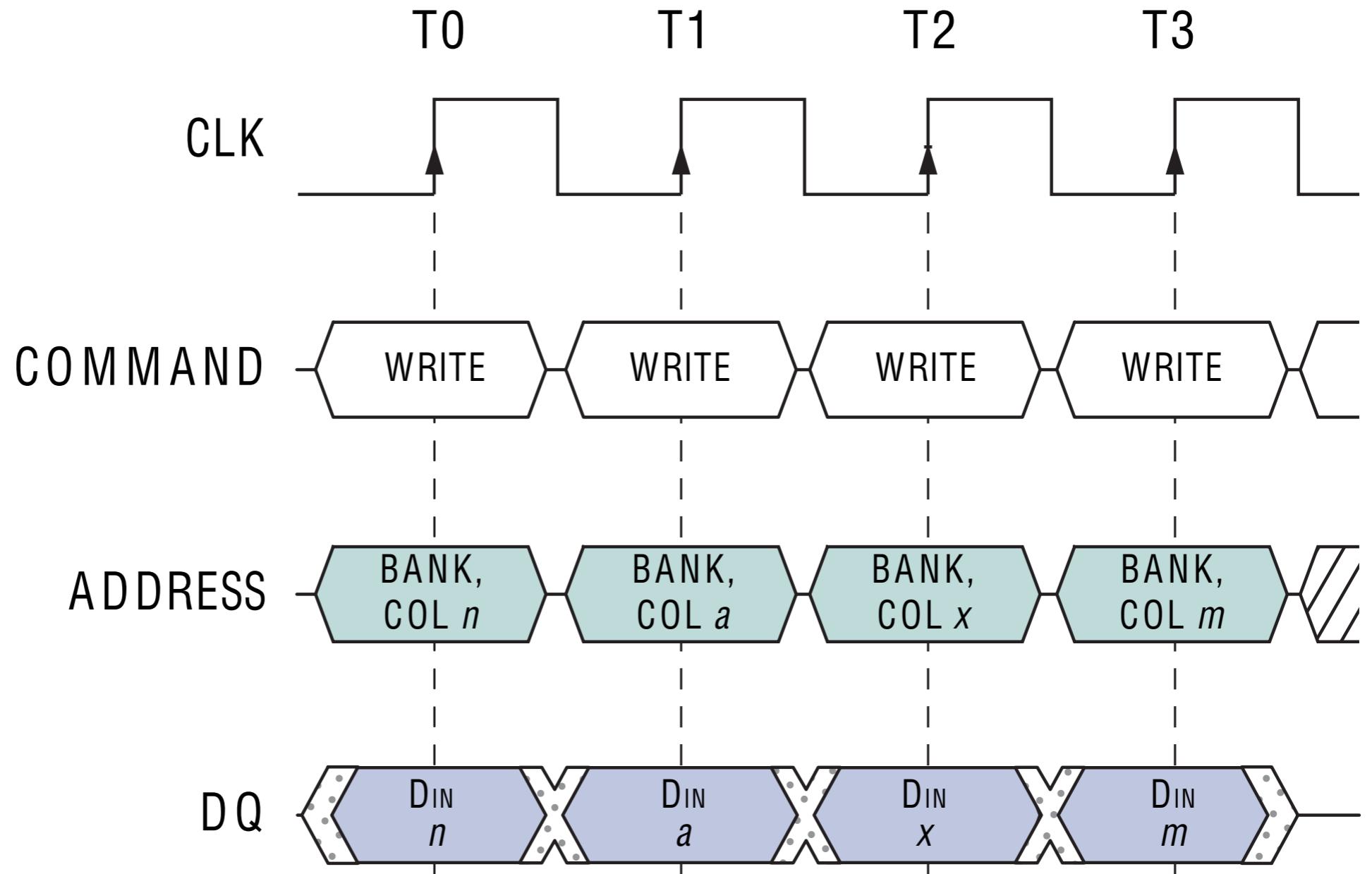


TRANSITIONING DATA

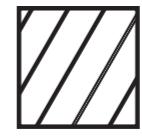


DON'T CARE

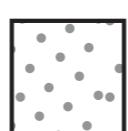
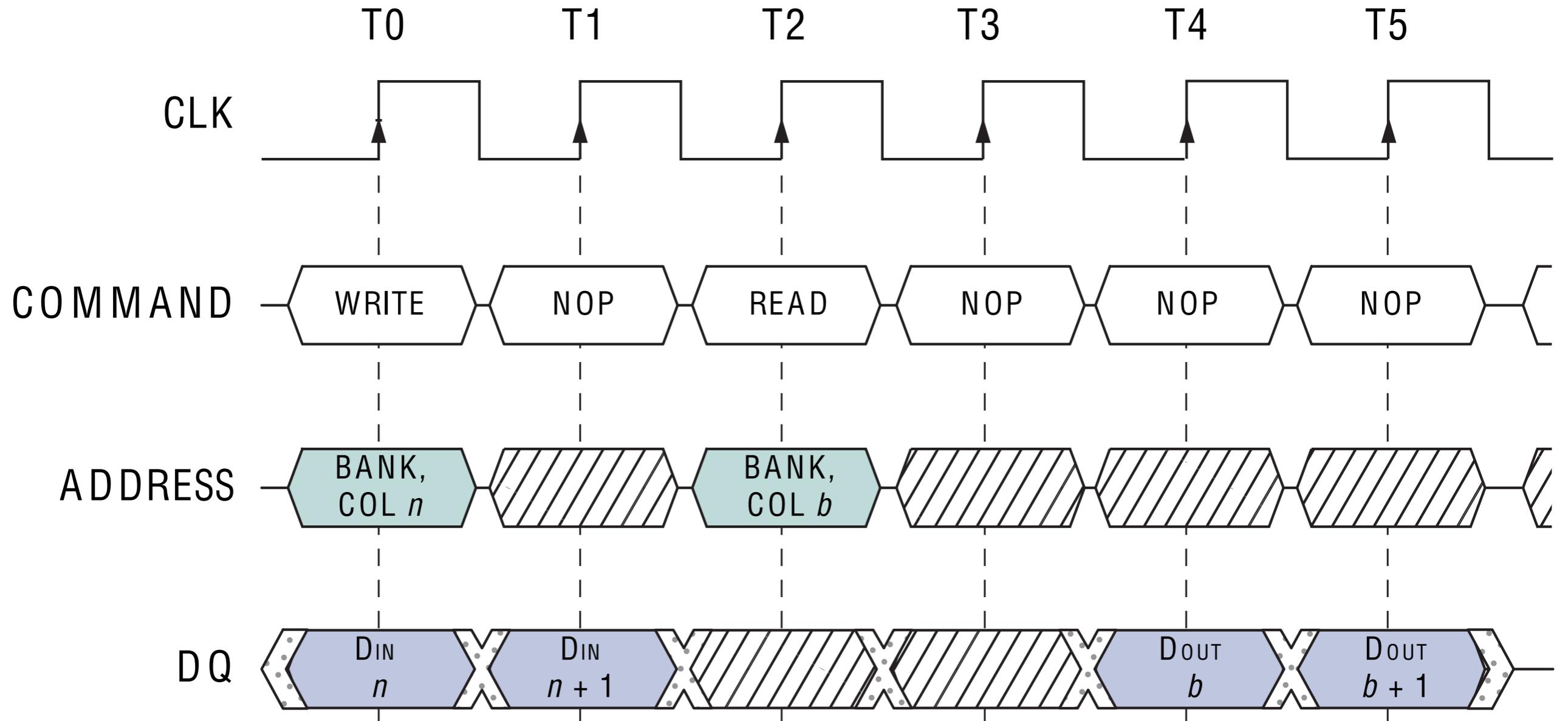
Random WRITE Cycles



Transitioning Data



Don't Care



Transitioning Data



Don't Care

