

Design of Digital Circuits

Lecture 22a: Memory Organization and Memory Technology

Prof. Onur Mutlu

ETH Zurich

Spring 2019

16 May 2019

Readings for This Lecture and Next

- Memory Hierarchy and Caches

- Required

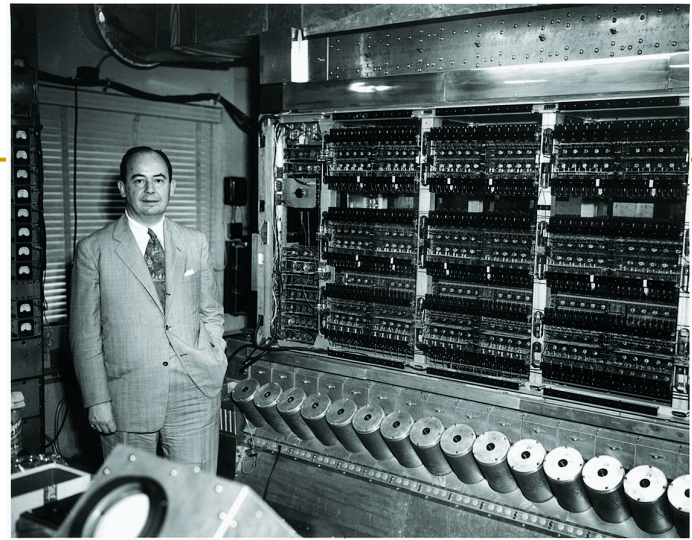
- H&H Chapters 8.1-8.3
- Refresh: P&P Chapter 3.5

- Recommended

- An early cache paper by Maurice Wilkes
 - Wilkes, “**Slave Memories and Dynamic Storage Allocation**,” IEEE Trans. On Electronic Computers, 1965.

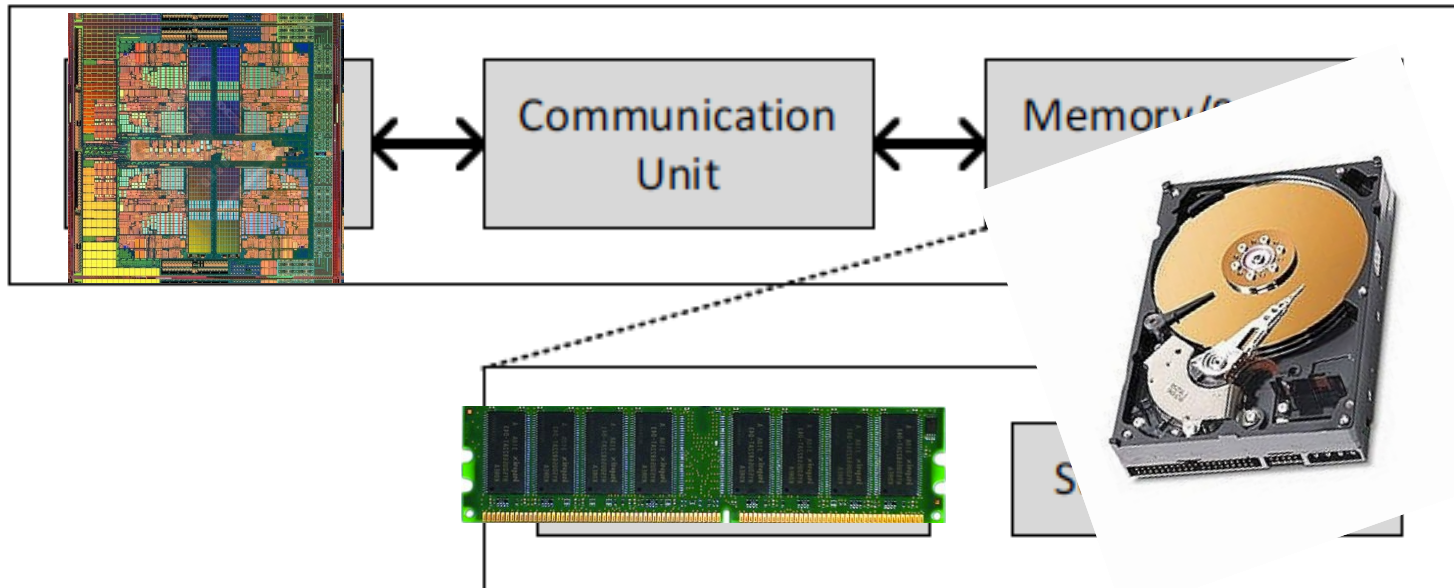
A Computing System

- Three key components
- Computation
- Communication
- Storage/memory



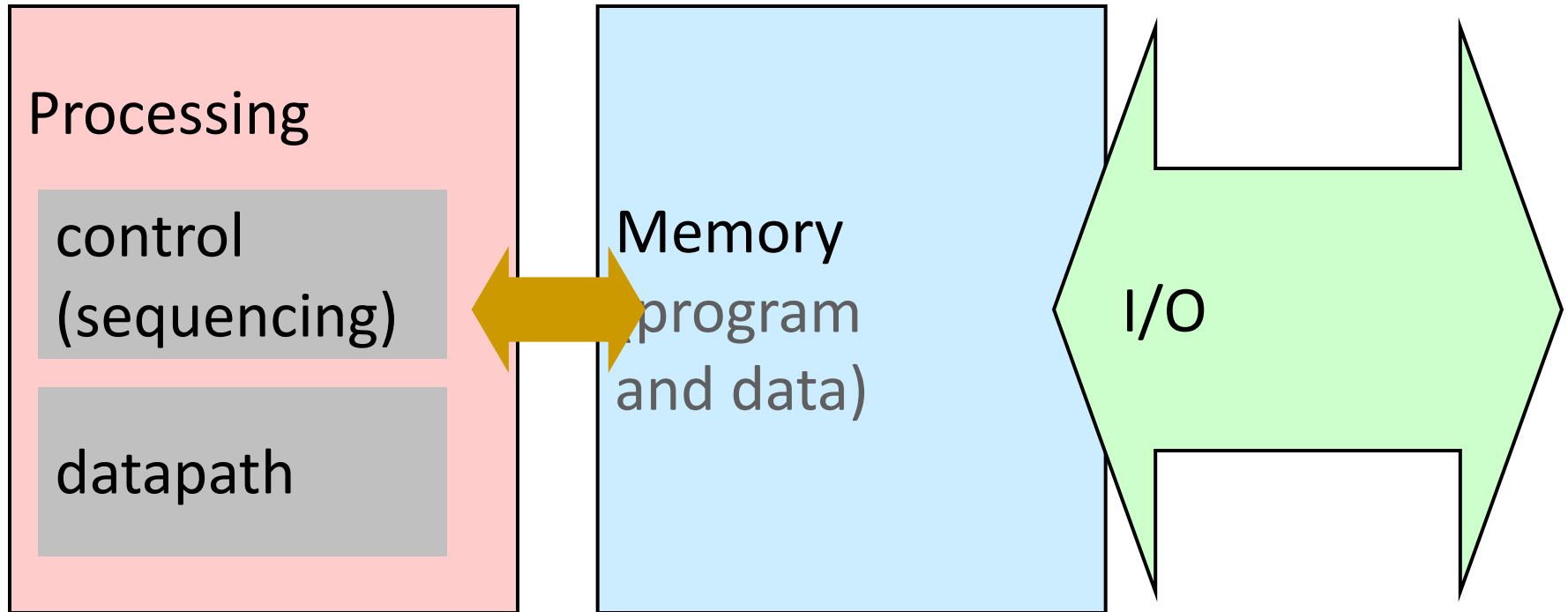
Burks, Goldstein, von Neumann, "Preliminary discussion of the logical design of an electronic computing instrument," 1946.

Computing System

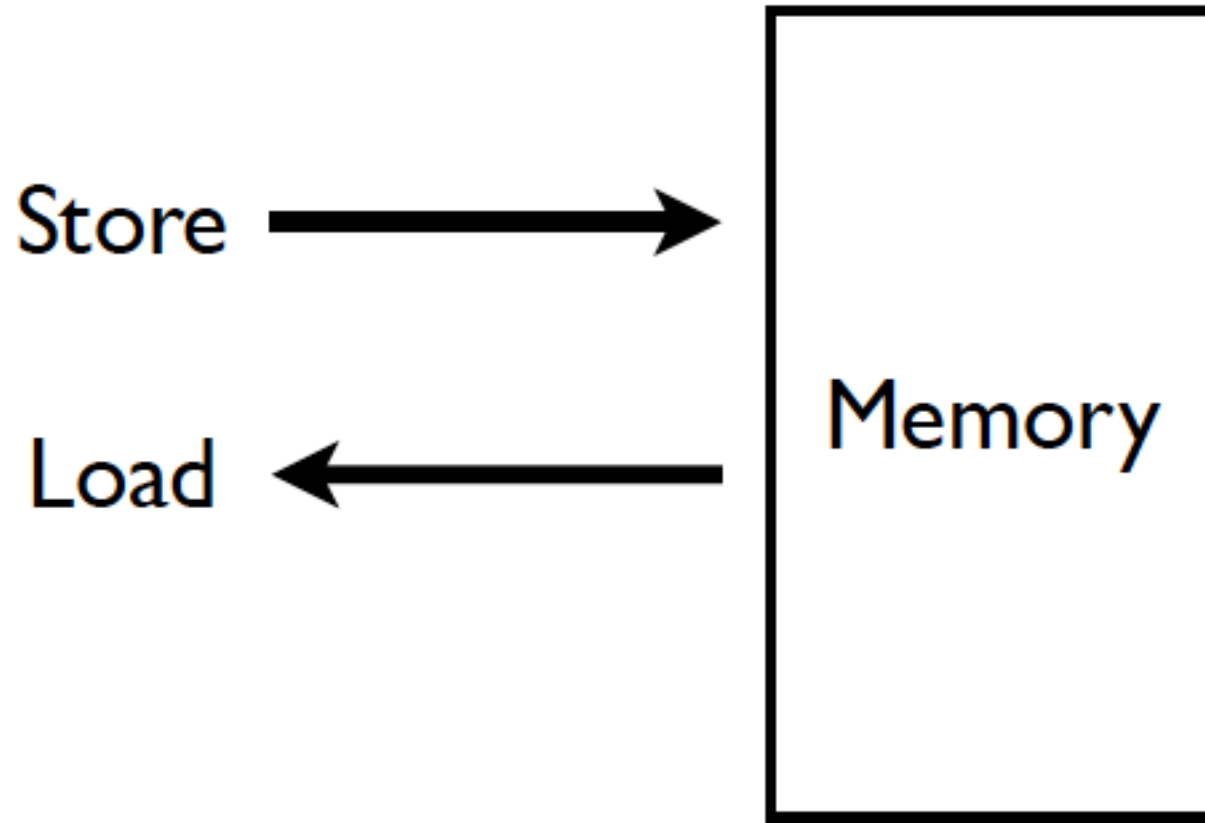


What is A Computer?

- We will cover all three components



Memory (Programmer's View)



Abstraction: Virtual vs. Physical Memory

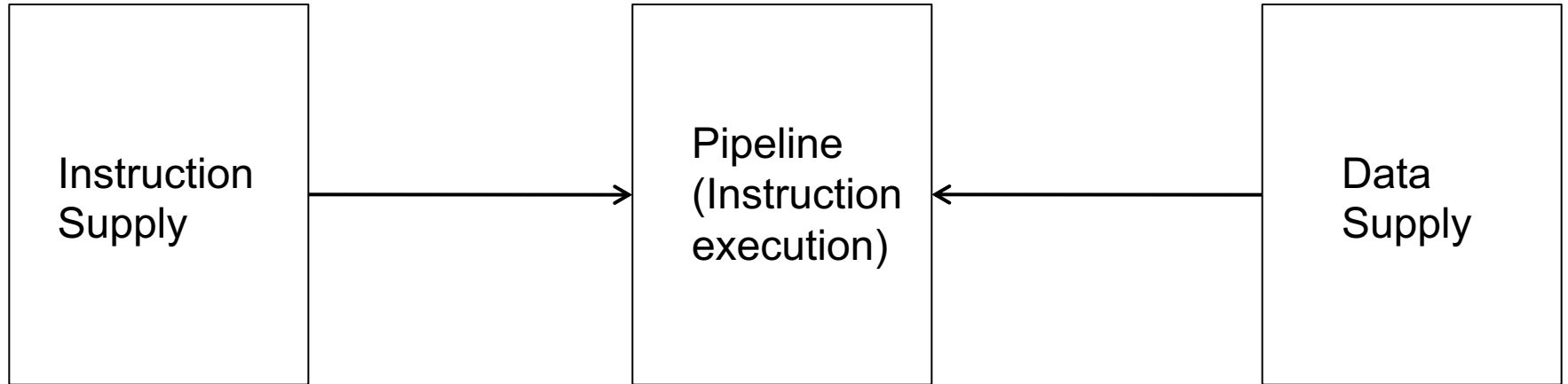
- **Programmer** sees **virtual memory**
 - Can assume the memory is “infinite”
 - Reality: **Physical memory** size is much smaller than what the programmer assumes
 - **The system** (system software + hardware, cooperatively) maps **virtual memory addresses** to **physical memory**
 - The system automatically manages the physical memory space **transparently to the programmer**
- + Programmer does not need to know the physical size of memory nor manage it → A small physical memory can appear as a huge one to the programmer → Life is easier for the programmer
- More complex system software and architecture

A classic example of the programmer/(micro)architect tradeoff

(Physical) Memory System

- You need a larger level of storage to manage a small amount of physical memory automatically
→ Physical memory has a backing store: disk
- We will first start with the physical memory system
- For now, ignore the virtual→physical indirection
- We will get back to it later, if time permits...

Idealism



- Zero latency access

- Infinite capacity

- Zero cost

- Perfect control flow

- No pipeline stalls

- Perfect data flow
(reg/memory dependencies)

- Zero-cycle interconnect
(operand communication)

- Enough functional units

- Zero latency compute

- Zero latency access

- Infinite capacity

- Infinite bandwidth

- Zero cost

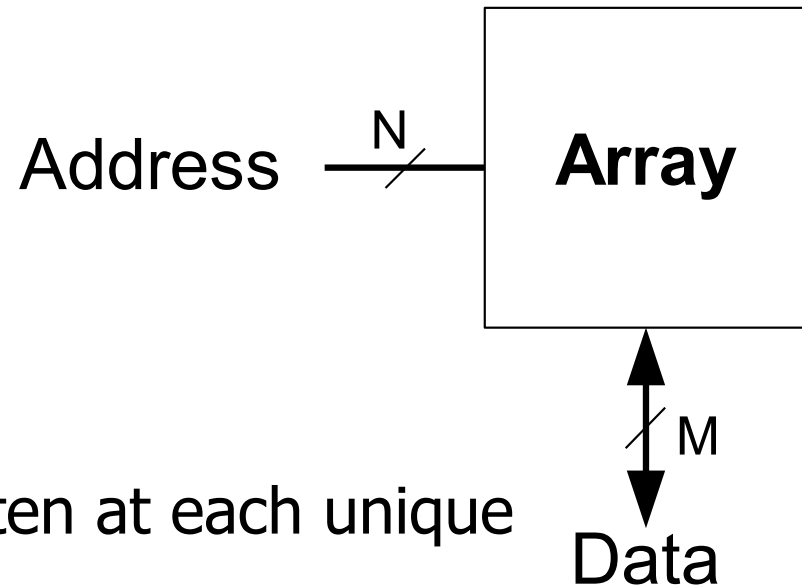
Quick Overview of Memory Arrays

How Can We Store Data?

- Flip-Flops (or Latches)
 - Very fast, parallel access
 - Very expensive (one bit costs tens of transistors)
 - Static RAM (we will describe them in a moment)
 - Relatively fast, only one data word at a time
 - Expensive (one bit costs 6 transistors)
 - Dynamic RAM (we will describe them a bit later)
 - Slower, one data word at a time, reading destroys content (refresh), needs special process for manufacturing
 - Cheap (one bit costs only one transistor plus one capacitor)
 - Other storage technology (flash memory, hard disk, tape)
 - Much slower, access takes a long time, non-volatile
 - Very cheap (no transistors directly involved)
-

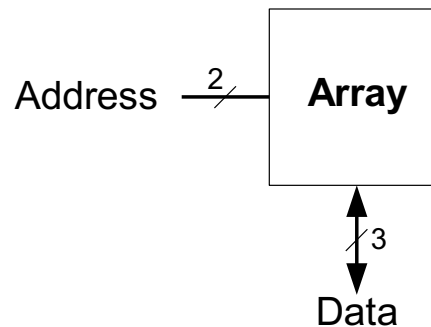
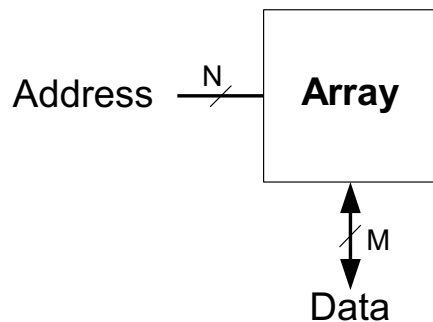
Array Organization of Memories

- 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

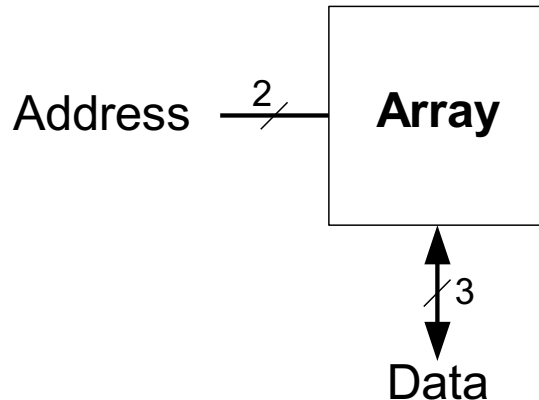
- Two-dimensional array of bit cells
 - Each bit cell stores one bit
- An array with N address bits and M data bits:
 - 2^N rows and M columns
 - Depth: number of rows (number of words)
 - Width: number of columns (size of word)
 - Array size: depth \times width = $2^N \times M$



Address	Data			
11	0	1	0	depth ↑ ↓
10	1	0	0	
01	1	1	0	
00	0	1	1	
	width ← →			

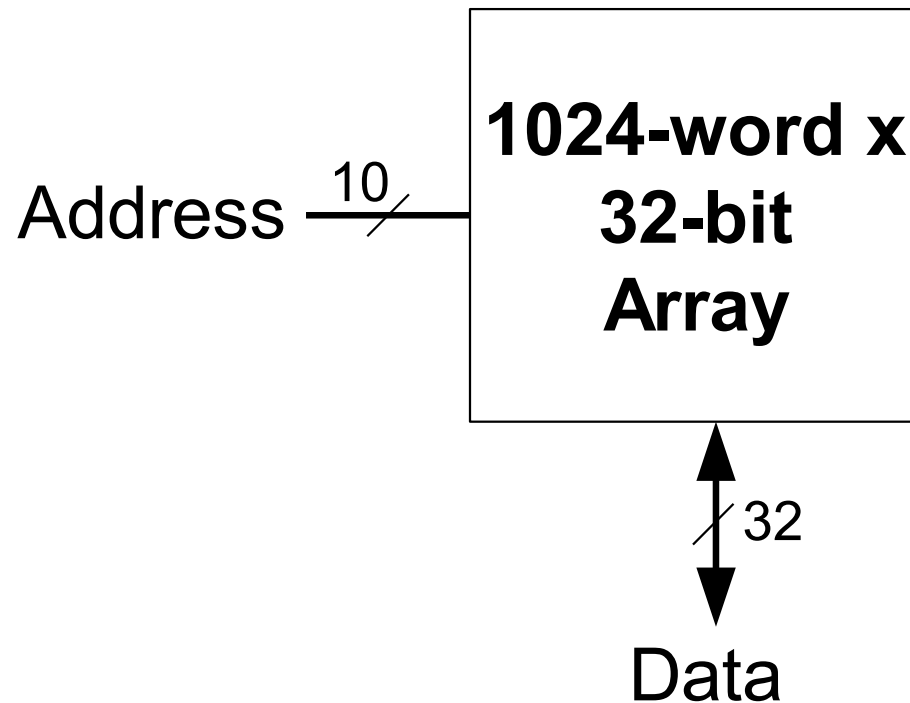
Memory Array Example

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



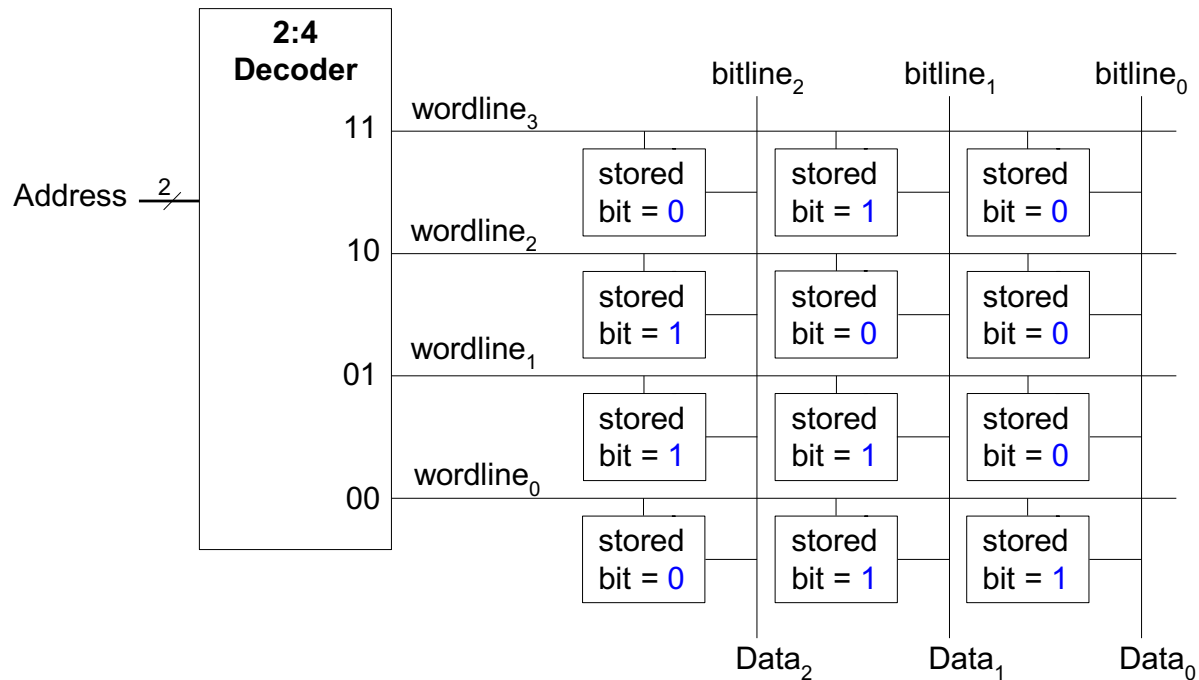
Address	Data			
11	0	1	0	depth ↑ ↓
10	1	0	0	
01	1	1	0	
00	0	1	1	
	width ←→			

Larger and Wider Memory Array Example



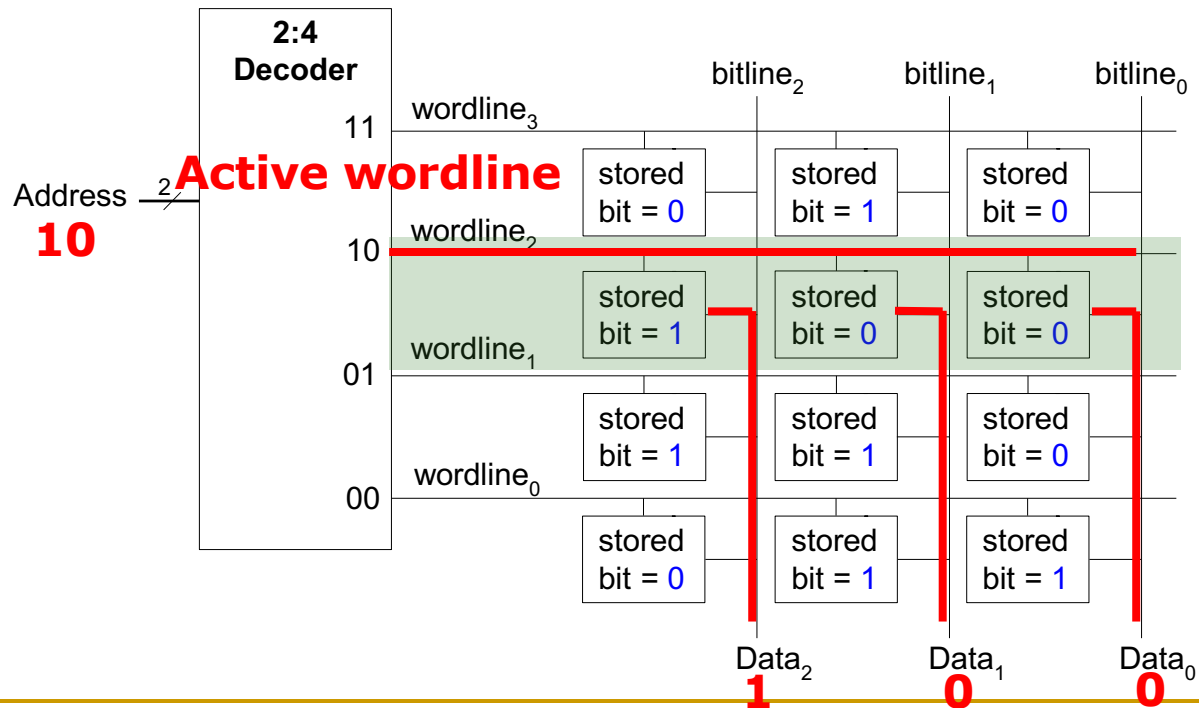
Memory Array Organization (I)

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



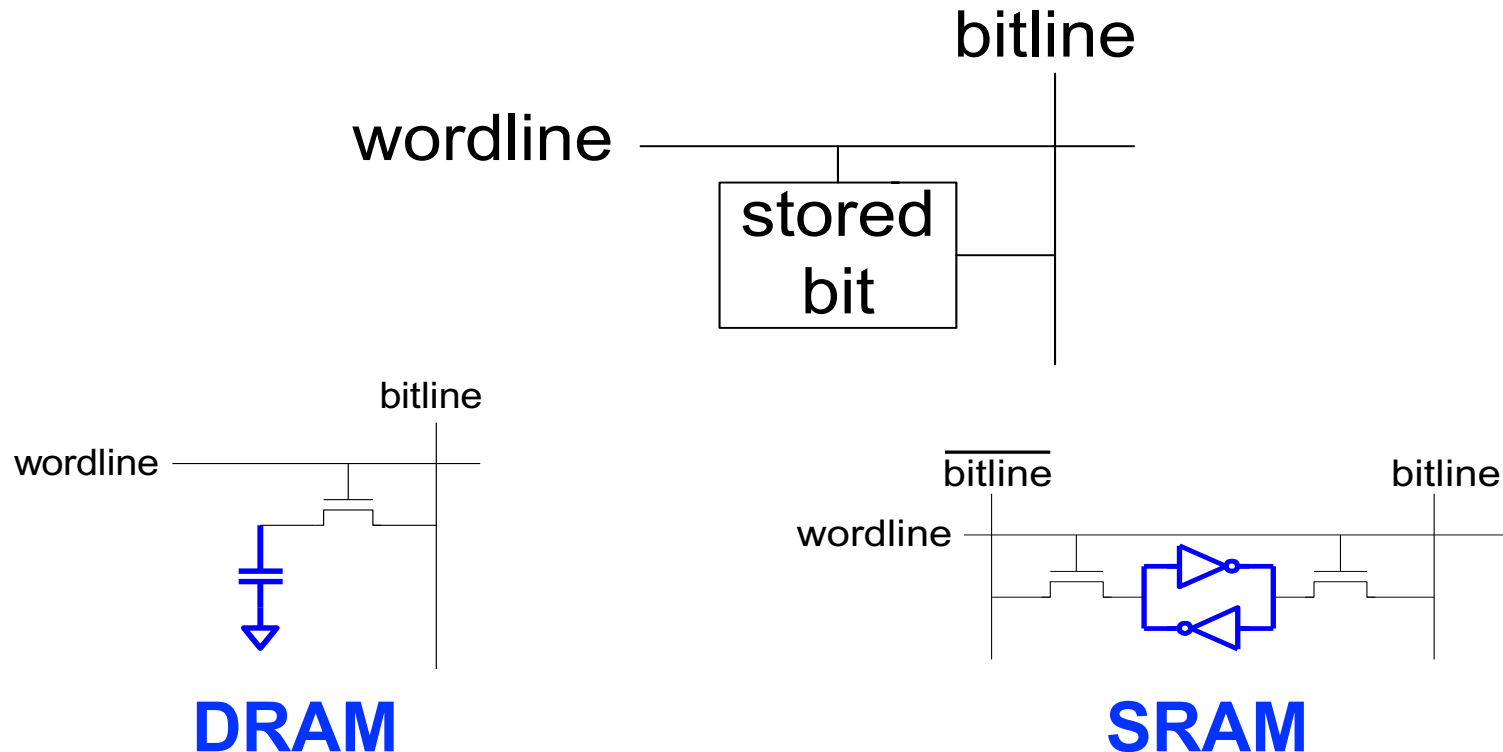
Memory Array Organization (II)

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



How is Access Controlled?

- Access transistors configured as switches connect the bit storage to the bitline.
- Access controlled by the wordline



Building Larger Memories

- Requires larger memory arrays
- Large → slow
- How do we make the memory large without making it very 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)

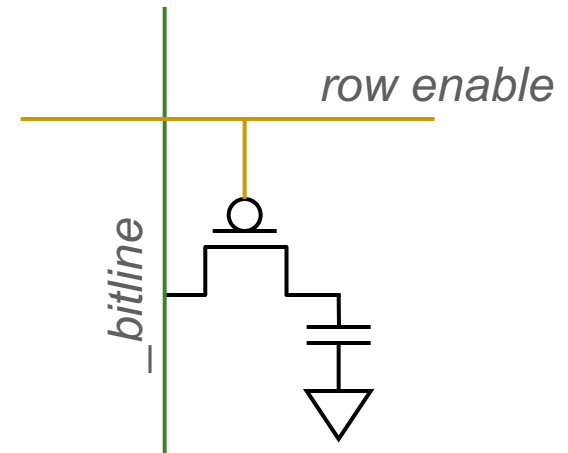
■ Interleaving (banking)

- ❑ **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?)

Memory Technology: DRAM and SRAM

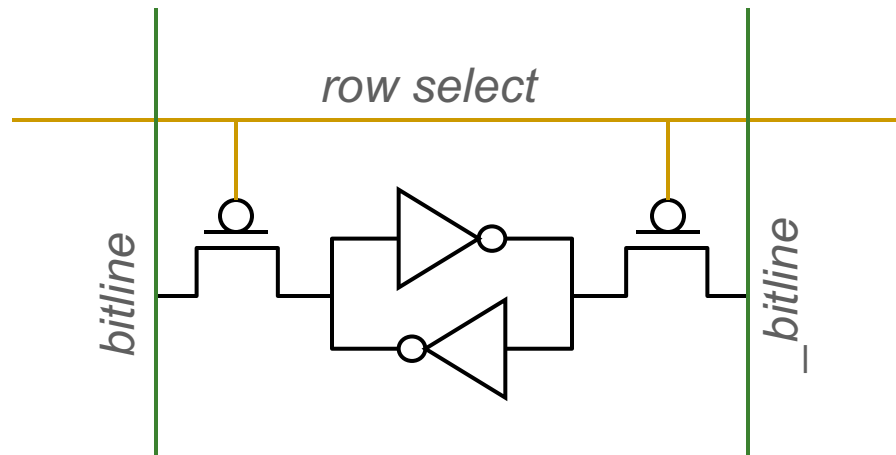
Memory Technology: DRAM

- Dynamic random access memory
- Capacitor charge state indicates stored value
 - Whether the capacitor is charged or discharged indicates storage of 1 or 0
 - 1 capacitor
 - 1 access transistor
- Capacitor leaks through the RC path
 - DRAM cell loses charge over time
 - DRAM cell needs to be refreshed

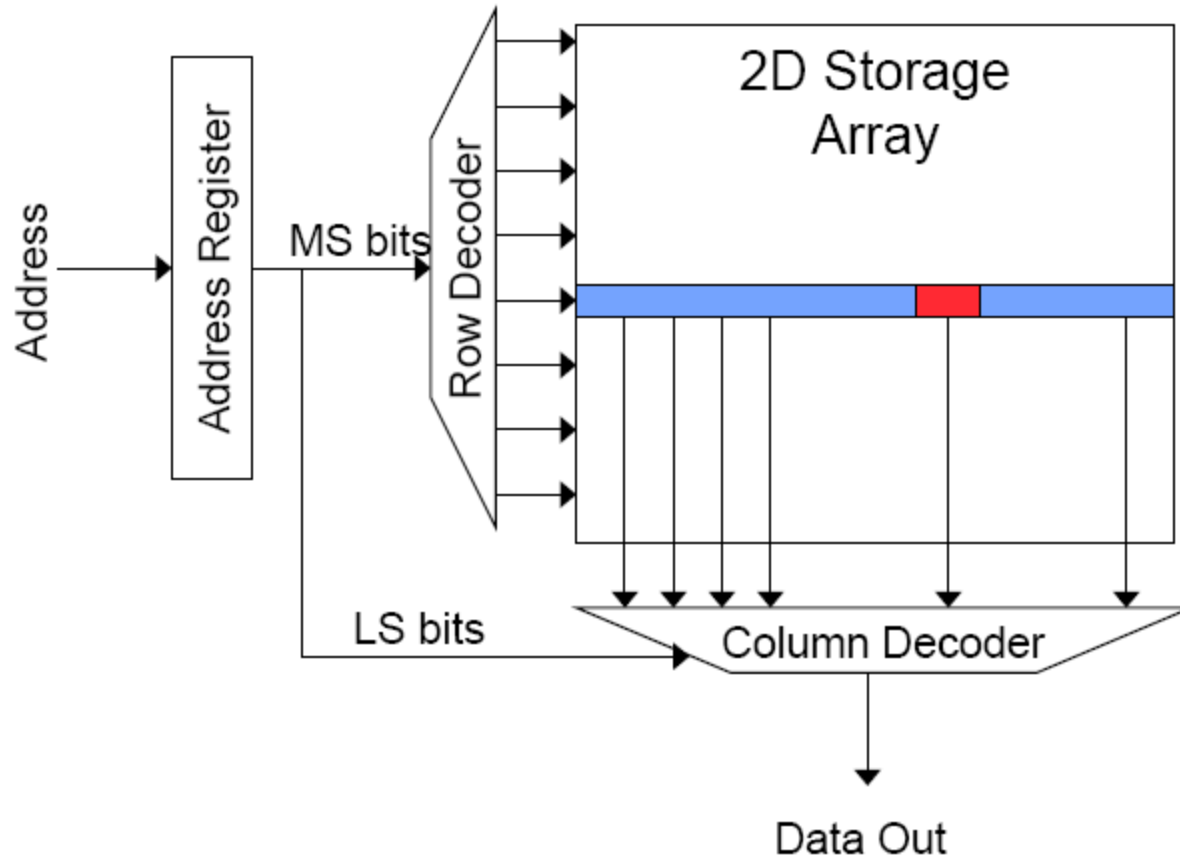


Memory Technology: SRAM

- Static random access memory
- Two cross coupled inverters store a single bit
 - Feedback path enables the stored value to persist in the “cell”
 - 4 transistors for storage
 - 2 transistors for access



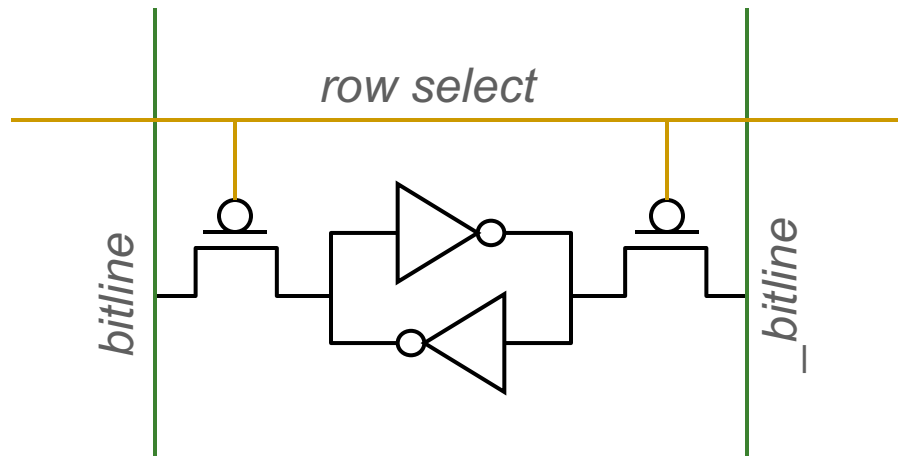
Memory Bank Organization and Operation



■ Read access sequence:

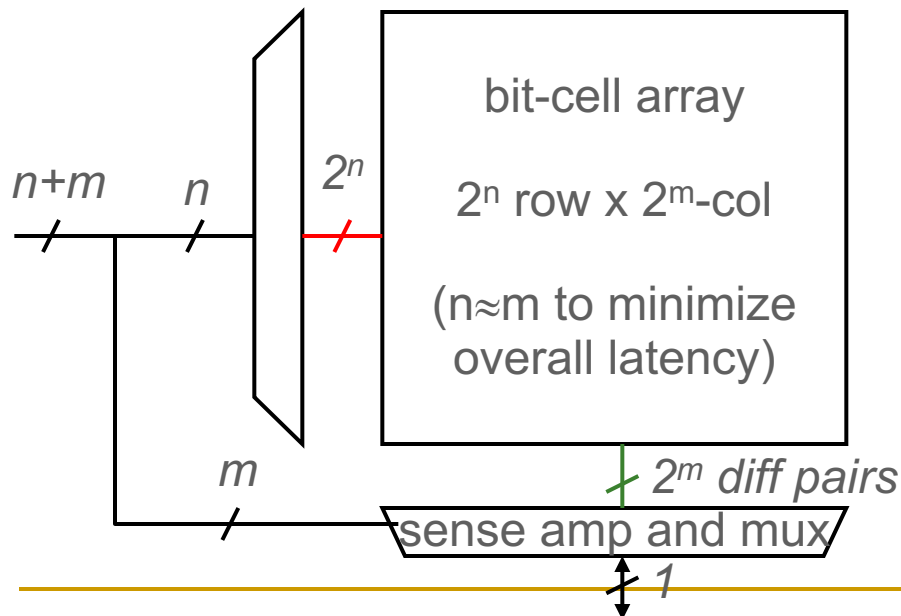
1. Decode row address & drive word-lines
2. Selected bits drive bit-lines
 - Entire row read
3. Amplify row data
4. Decode column address & select subset of row
 - Send to output
5. Precharge bit-lines
 - For next access

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)

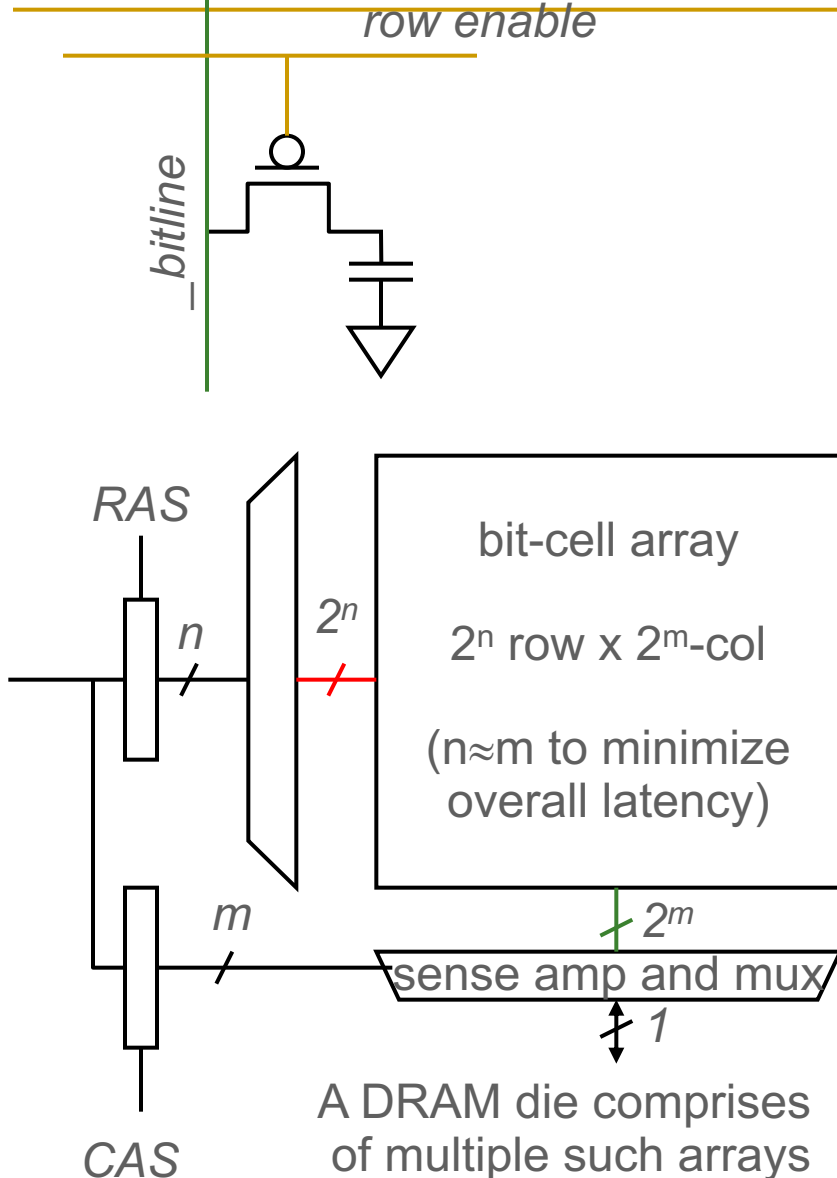


Access latency dominated by steps 2 and 3

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)



Bits stored as charges on node capacitance (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

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)

Design of Digital Circuits

Lecture 22a: Memory Organization and Memory Technology

Prof. Onur Mutlu

ETH Zurich

Spring 2019

16 May 2019