

The Memory Hierarchy

CSCI 2021: Machine Architecture and Organization

Antonia Zhai

Department Computer Science and Engineering

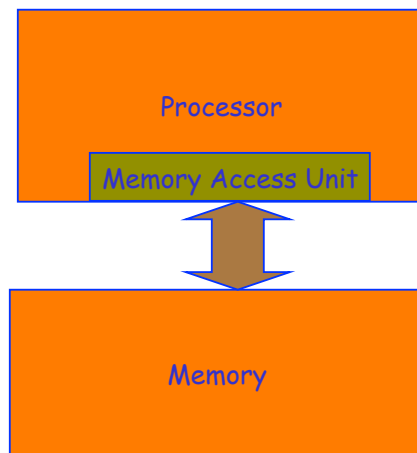
University of Minnesota

<http://www.cs.umn.edu/~zhai>

With Slides from Bryant and O'Hallaron



Abstraction



**RAM: Random
Access Memory**

Is this abstraction accurate? How does the memory really work?

Random-Access Memory (RAM)

Key features

- **RAM** is packaged as a chip.
- Basic storage unit is a **cell** (one bit per cell).
- Multiple RAM chips form a memory.

Static RAM (**SRAM**)

- Each cell stores bit with a six-transistor circuit.
- Retains value indefinitely, as long as it is kept powered.
- Relatively insensitive to disturbances such as electrical noise.
- Faster and more expensive than DRAM.

Dynamic RAM (**DRAM**)

- Each cell stores bit with a capacitor and transistor.
- Value must be refreshed every 10-100 ms.
- Sensitive to disturbances.
- Slower and cheaper than SRAM.

3/25/15

CSCI 2021

3

SRAM vs DRAM Summary

	Tran. per bit	Access time	Persist?	Sensitive?	Cost	Applications
SRAM	6	1X	Yes	No	100x	cache memories
DRAM	1	10X	No	Yes	1X	Main memories, frame buffers

3/25/15

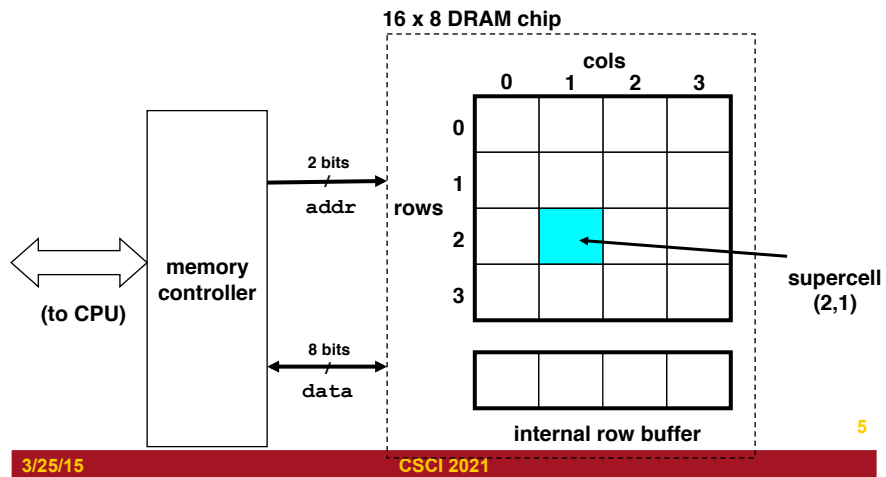
CSCI 2021

4

Conventional DRAM Organization

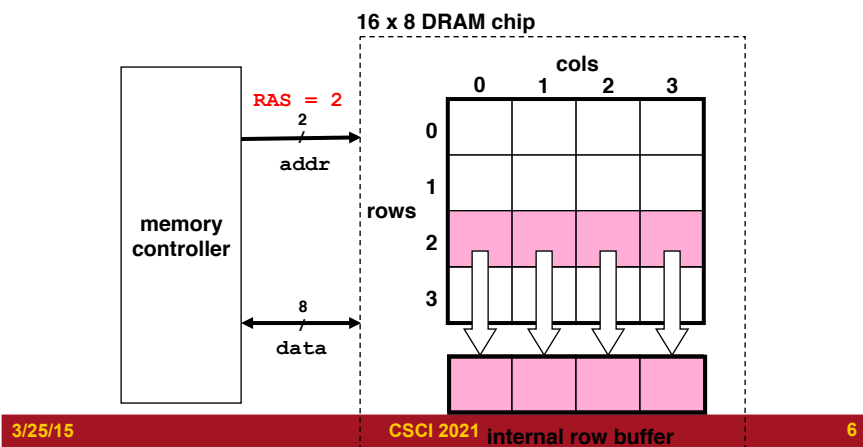
$d \times w$ DRAM:

- dw total bits organized as d **supercells** of size w bits



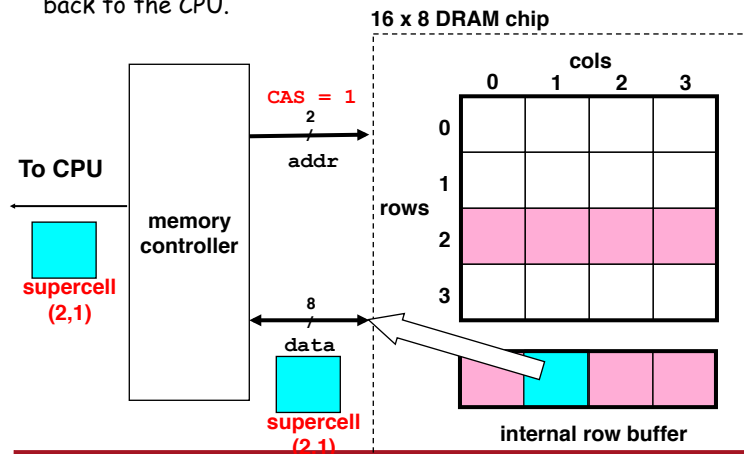
Reading DRAM Supercell (2,1)

- Step 1(a): Row access strobe (**RAS**) selects row 2.
- Step 1(b): Row 2 copied from DRAM array to row buffer.



Reading DRAM Supercell (2,1)

- Step 2(a): Column access strobe (**CAS**) selects column 1.
- Step 2(b): Supercell (2,1) copied from buffer to data lines, and eventually back to the CPU.

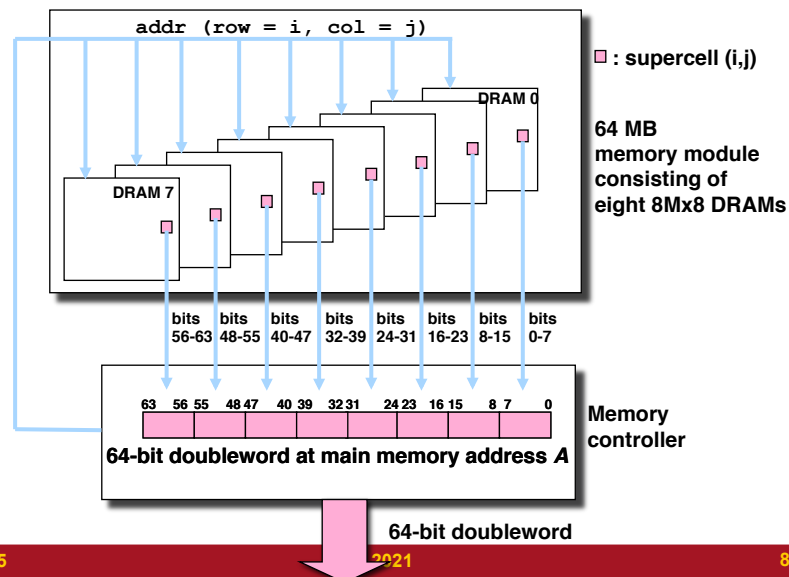


3/25/15

CSCI 2021

7

Memory Modules



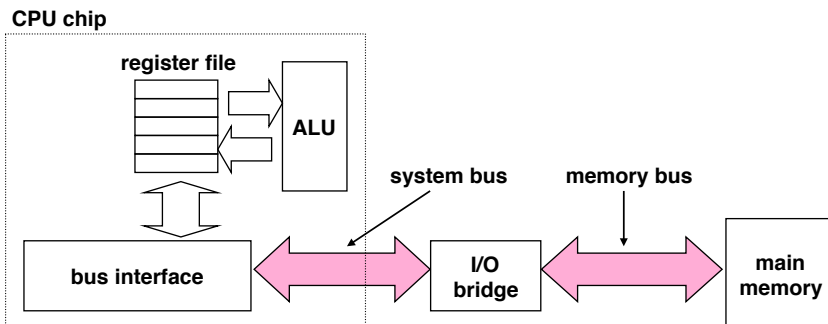
3/25/15

2021

8

Bus Structure Connecting CPU and Memory

- A **bus** is a collection of parallel wires that carry address, data, and control signals.
- Buses are typically shared by multiple devices.



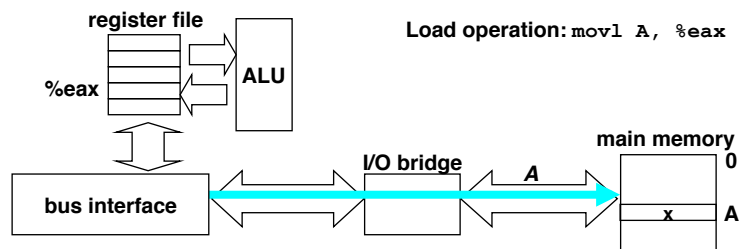
3/25/15

CSCI 2021

9

Memory Read Transaction (1)

- CPU places address *A* on the memory bus.



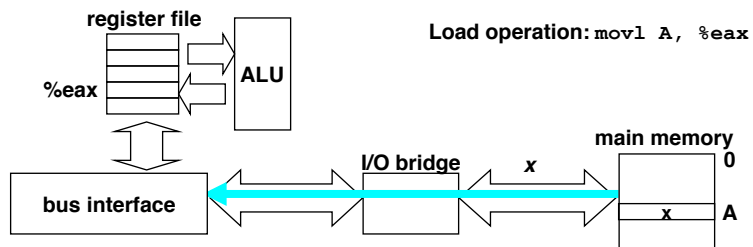
3/25/15

CSCI 2021

10

Memory Read Transaction (2)

- Main memory reads *A* from the memory bus, retrieves word *x*, and places it on the bus.



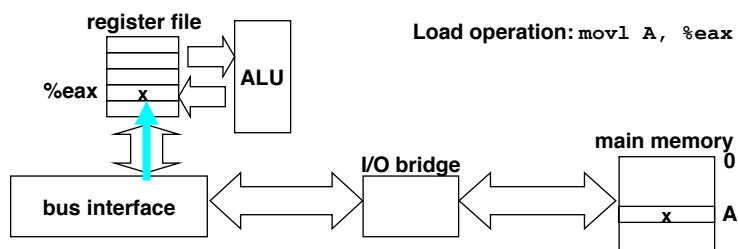
3/25/15

CSCI 2021

11

Memory Read Transaction (3)

- CPU read word *x* from the bus and copies it into register *%eax*.



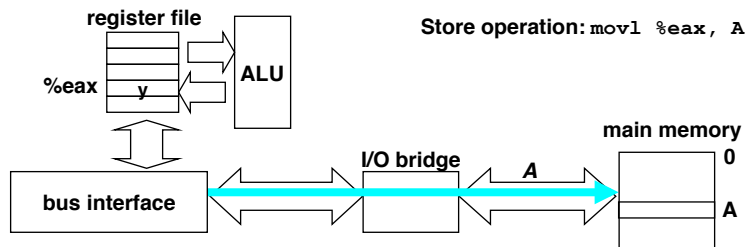
3/25/15

CSCI 2021

12

Memory Write Transaction (1)

- CPU places address *A* on bus. Main memory reads it and waits for the corresponding data word to arrive.



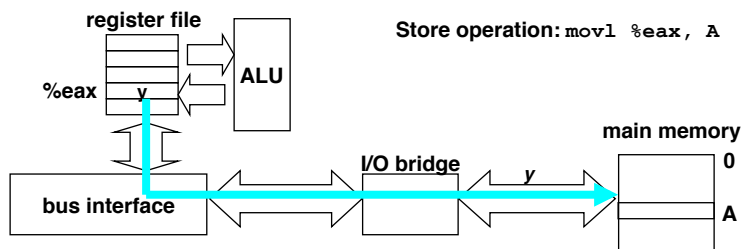
3/25/15

CSCI 2021

13

Memory Write Transaction (2)

- CPU places data word *y* on the bus.



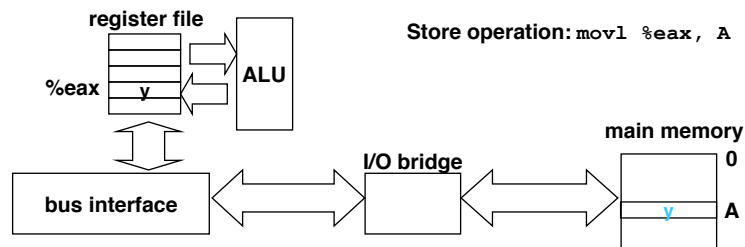
3/25/15

CSCI 2021

14

Memory Write Transaction (3)

- Main memory read data word y from the bus and stores it at address A .



3/25/15

CSCI 2021

15

Enhanced DRAMs

All enhanced DRAMs are built around the conventional DRAM core.

- Fast page mode DRAM (**FPM DRAM**)
 - Access contents of row with $[RAS, CAS, CAS, CAS, CAS]$ instead of $[(RAS, CAS), (RAS, CAS), (RAS, CAS), (RAS, CAS)]$.
- Extended data out DRAM (**EDO DRAM**)
 - Enhanced FPM DRAM with more closely spaced CAS signals.
- Synchronous DRAM (**SDRAM**)
 - Driven with rising clock edge instead of asynchronous control signals.
- Double data-rate synchronous DRAM (**DDR SDRAM**) and (**DDR2 SDRAM**)
 - Enhancement of SDRAM that uses both clock edges as control signals.
 - DDR2 clocks the bus at twice the speed of the memory cells
- Video RAM (**VRAM**)
 - Like FPM DRAM, but output is produced by shifting row buffer
 - Dual ported (allows concurrent reads and writes)

3/25/15

CSCI 2021

16

Nonvolatile Memories

- **DRAM and SRAM are volatile memories**
 - Lose information if powered off.
- **Nonvolatile memories retain value even if powered off**
 - Read-only memory (**ROM**): programmed during production
 - Programmable ROM (**PROM**): can be programmed once
 - Erasable PROM (**EPROM**): can be bulk erased (UV, X-Ray)
 - Electrically erasable PROM (**EEPROM**): electronic erase capability
 - Flash memory: EEPROMs with partial (sector) erase capability
 - Wears out after about 100,000 erasings.
- **Uses for Nonvolatile Memories**
 - Firmware programs stored in a ROM (BIOS, controllers for disks, network cards, graphics accelerators, security subsystems,...)
 - Solid state disks (replace rotating disks in thumb drives, smart phones, mp3 players, tablets, laptops,...)
 - Disk caches

3/25/15

CSCI 2021

17

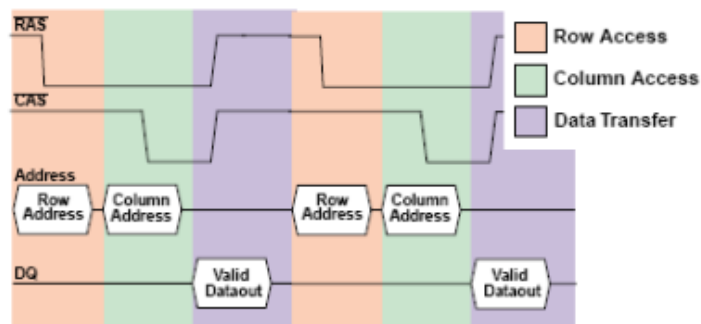
DRAM TUTORIAL

ISCA 2002

Bruce Jacob
David Wang

University of
Maryland

Read Timing for Conventional DRAM

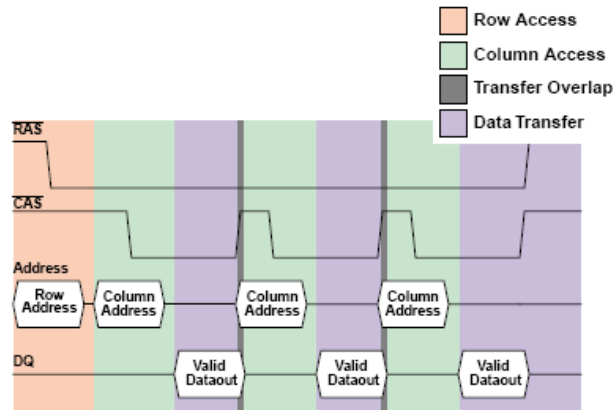


3/25/15

CSCI 2021

18

Read Timing for Fast Page Mode

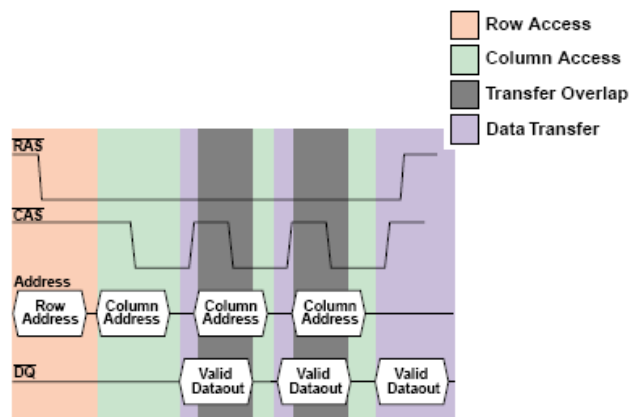


3/25/15

CSCI 2021

19

Read Timing for Extended Data Out

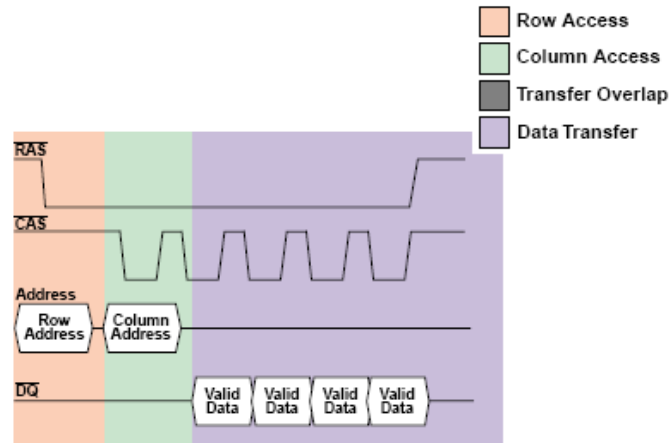


3/25/15

CSCI 2021

20

Read Timing for Pipeline Burst EDO

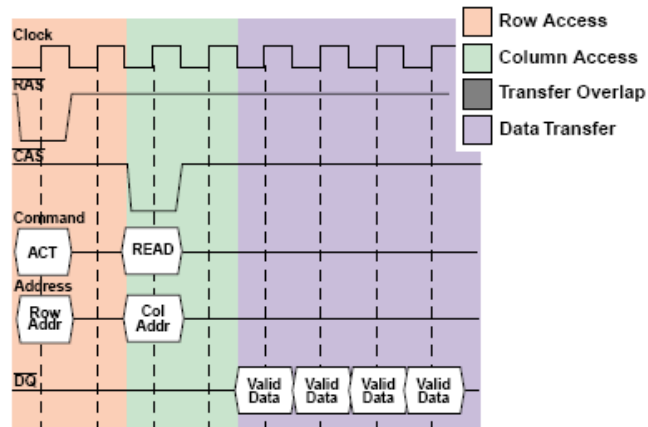


3/25/15

CSCI 2021

21

Read Timing for Synchronous DRAM

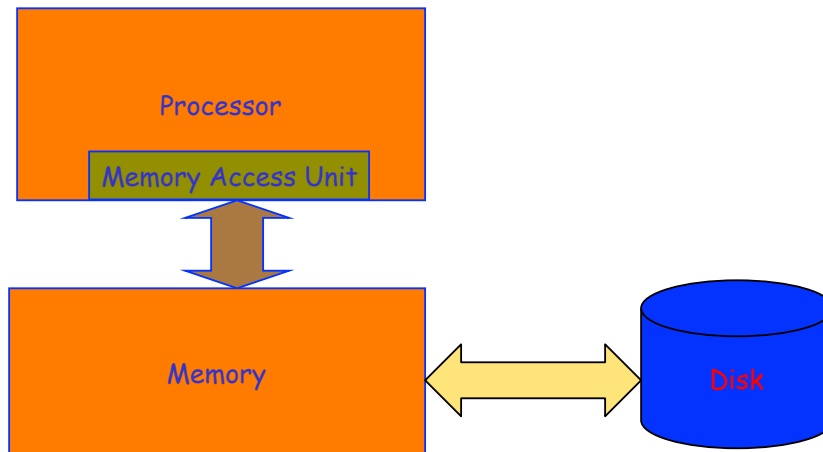


3/25/15

CSCI 2021

22

Our Abstraction



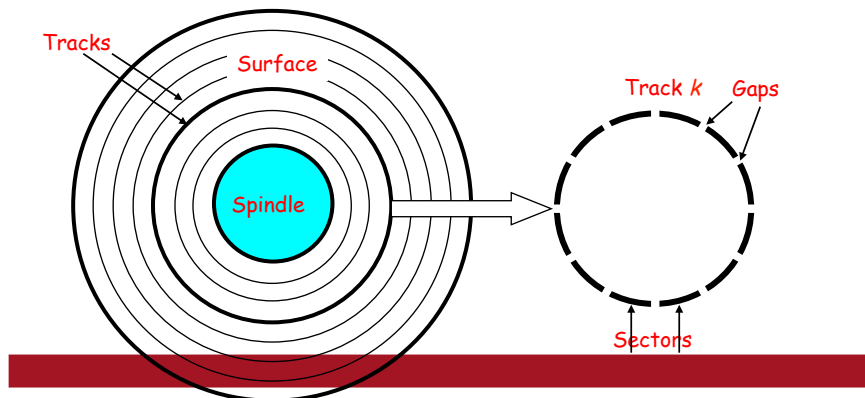
3/25/15

CSCI 2021

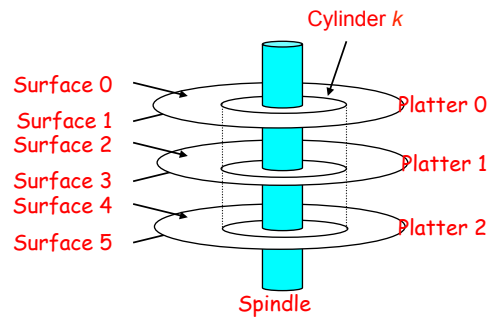
23

Disk Geometry

- Disks consist of **platters**, each with two **surfaces**.
- Each surface consists of concentric rings called **tracks**.
- Each track consists of **sectors** separated by **gaps**.



Disk Geometry (Multiple-Platter View)



- Aligned tracks form a cylinder.
-

Disk Capacity

- **Capacity**: maximum number of bits that can be stored.
 - Vendors express capacity in units of gigabytes (GB), where $1 \text{ GB} = 10^9 \text{ Bytes}$ (Lawsuit pending! Claims deceptive advertising).
 - Capacity is determined by these technology factors:
 - **Recording density** (bits/in): number of bits that can be squeezed into a 1 inch segment of a track.
 - **Track density** (tracks/in): number of tracks that can be squeezed into a 1 inch radial segment.
 - **Areal density** (bits/in²): product of recording and track density.
 - Modern disks partition tracks into disjoint subsets called **recording zones**
 - Each track in a zone has the same number of sectors, determined by the circumference of innermost track.
 - Each zone has a different number of sectors/track
-

Computing Disk Capacity

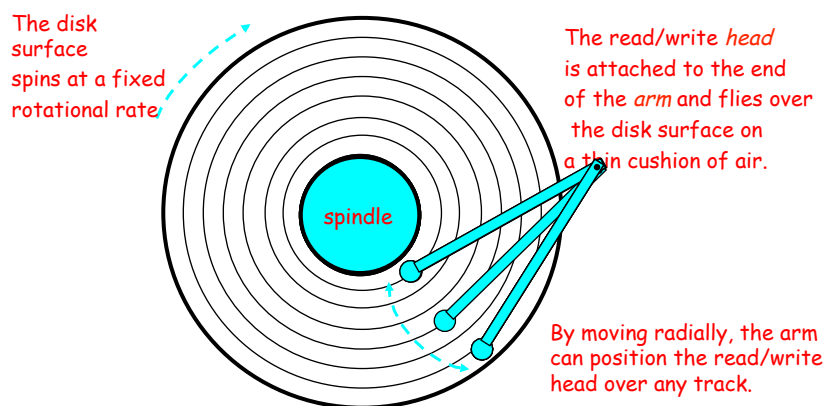
$$\text{Capacity} = (\# \text{ bytes/sector}) \times (\text{avg. } \# \text{ sectors/track}) \times$$
$$(\# \text{ tracks/surface}) \times (\# \text{ surfaces/platter}) \times$$
$$(\# \text{ platters/disk})$$

Example:

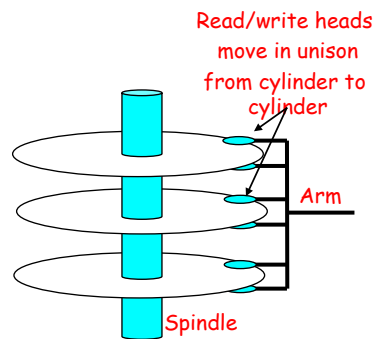
- 512 bytes/sector
- 300 sectors/track (on average)
- 20,000 tracks/surface
- 2 surfaces/platter
- 5 platters/disk

$$\begin{aligned}\text{Capacity} &= 512 \times 300 \times 20000 \times 2 \times 5 \\ &= 30,720,000,000 \\ &= 30.72 \text{ GB}\end{aligned}$$

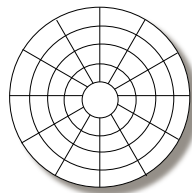
Disk Operation (Single-Platter View)



Disk Operation (Multi-Platter View)



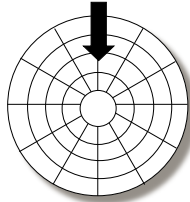
Disk Structure - top view of single platter



Surface organized into tracks

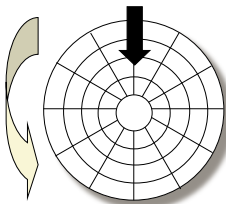
Tracks divided into sectors

Disk Access



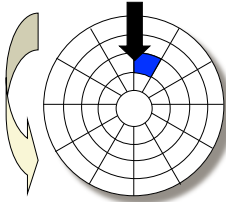
Head in position above a track

Disk Access



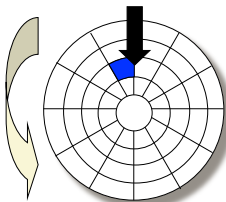
Rotation is counter-clockwise

Disk Access - Read



About to read blue sector

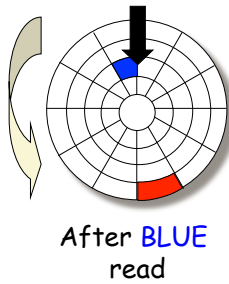
Disk Access - Read



After BLUE
read

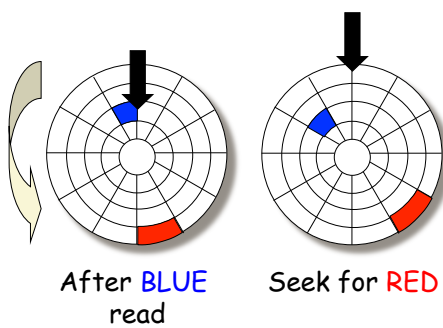
After reading blue sector

Disk Access - Read



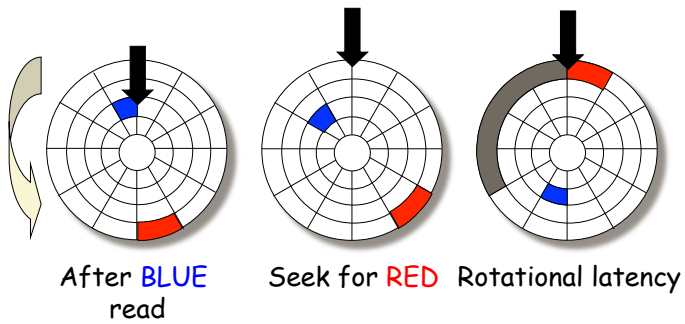
Red request scheduled next

Disk Access - Seek



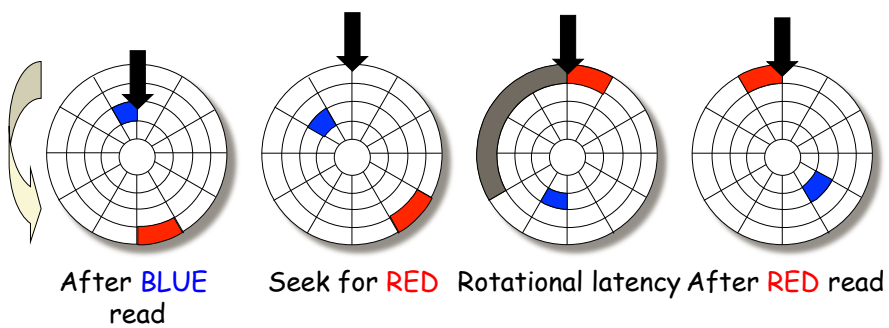
Seek to red's track

Disk Access - Rotational Latency



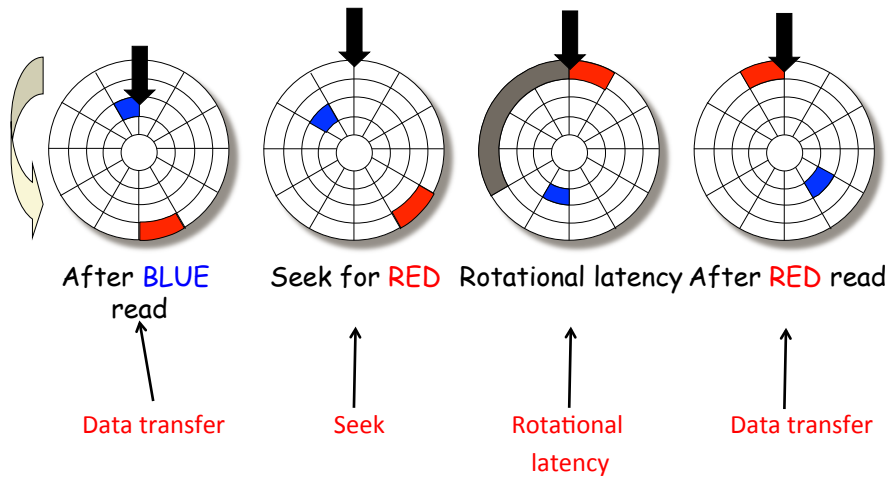
Wait for red sector to rotate around

Disk Access - Read



Complete read of red

Disk Access - Service Time Components



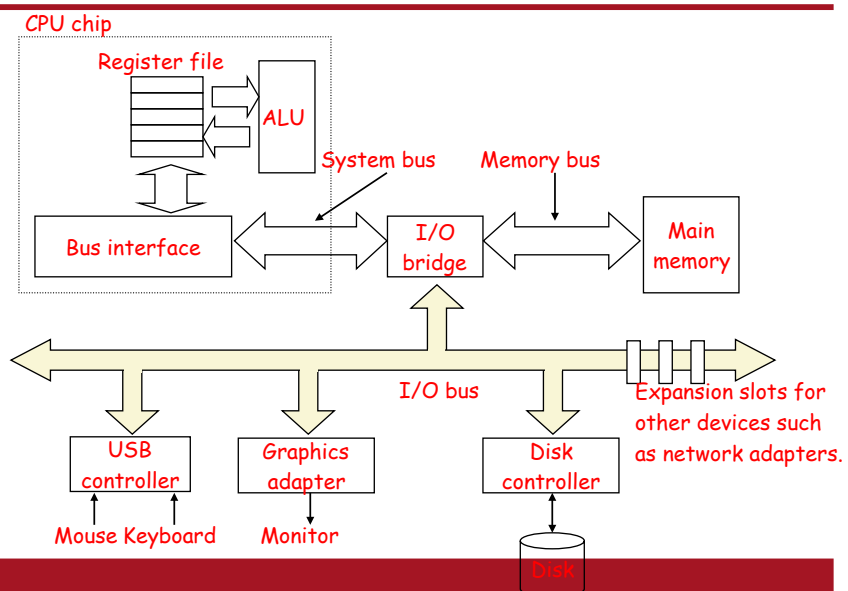
Disk Access Time Example

- Given:
 - Rotational rate = 7,200 RPM
 - Average seek time = 9 ms.
 - Avg # sectors/track = 400.
- Derived:
 - Tavg rotation = $\frac{1}{2} \times (60 \text{ secs} / 7200 \text{ RPM}) \times 1000 \text{ ms/sec} = 4 \text{ ms}$.
 - Tavg transfer = $60 / 7200 \text{ RPM} \times 1 / 400 \text{ secs/track} \times 1000 \text{ ms/sec} = 0.02 \text{ ms}$
 - Taccess = 9 ms + 4 ms + 0.02 ms
- Important points:
 - Access time dominated by seek time and rotational latency.
 - First bit in a sector is the most expensive, the rest are free.
 - SRAM access time is about 4 ns/doubleword, DRAM about 60 ns
 - Disk is about 40,000 times slower than SRAM,
 - 2,500 times slower than DRAM.

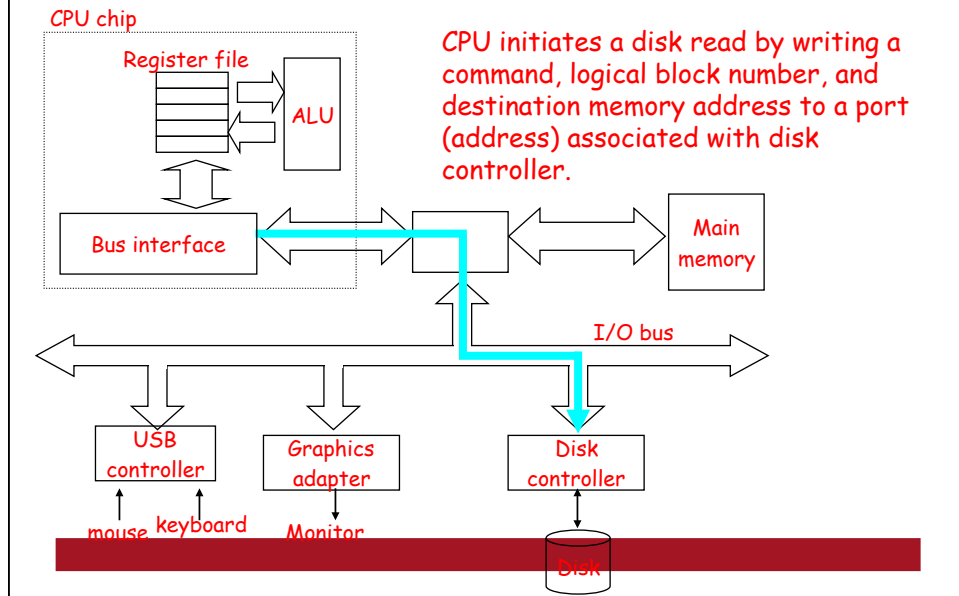
Logical Disk Blocks

- Modern disks present a simpler abstract view of the complex sector geometry:
 - The set of available sectors is modeled as a sequence of b-sized **logical blocks** (0, 1, 2, ...)
- Mapping between logical blocks and actual (physical) sectors
 - Maintained by hardware/firmware device called disk controller.
 - Converts requests for logical blocks into (surface, track, sector) triples.
- Allows controller to set aside spare cylinders for each zone.
 - Accounts for the difference in "formatted capacity" and "maximum capacity".

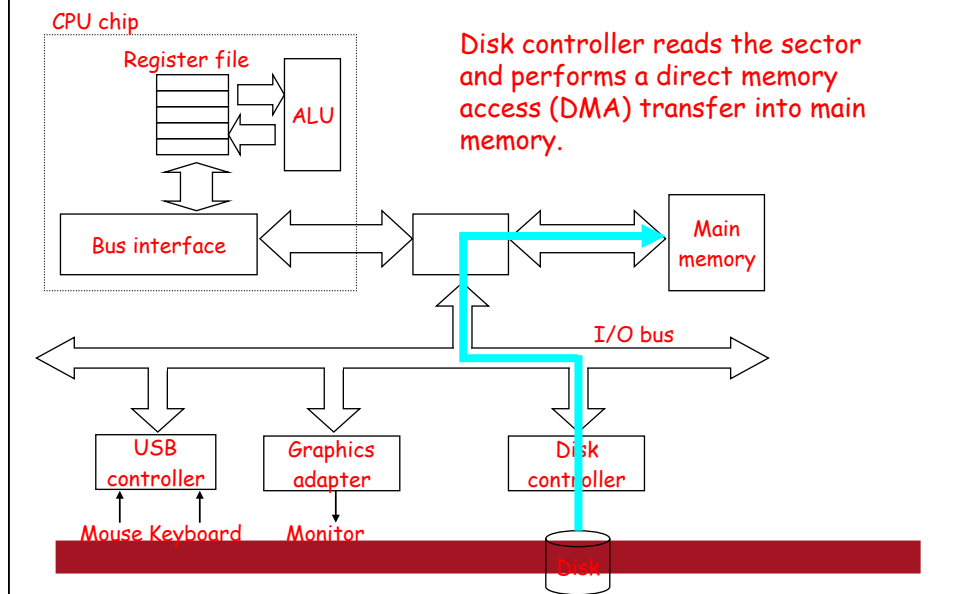
I/O Bus



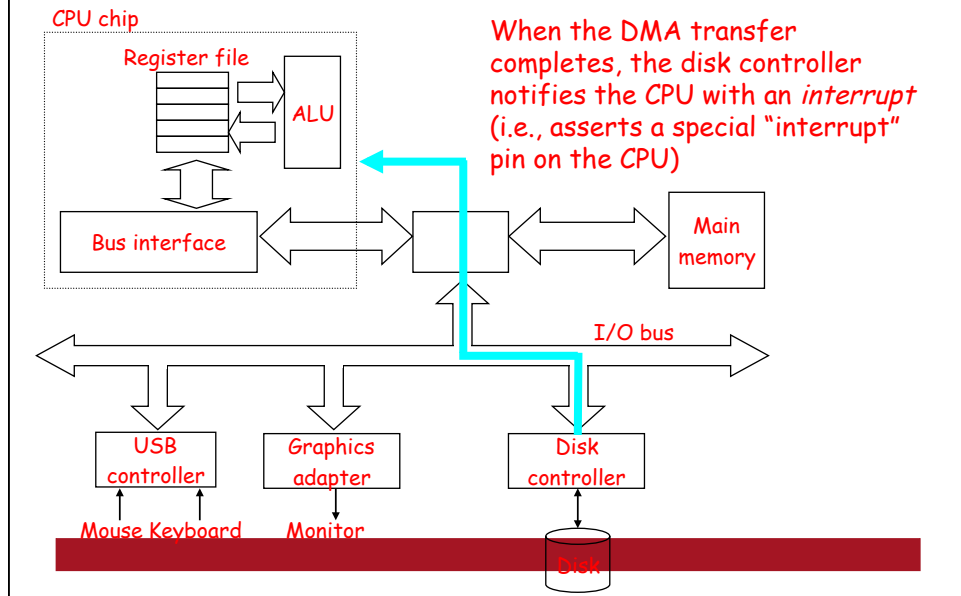
Reading a Disk Sector (1)



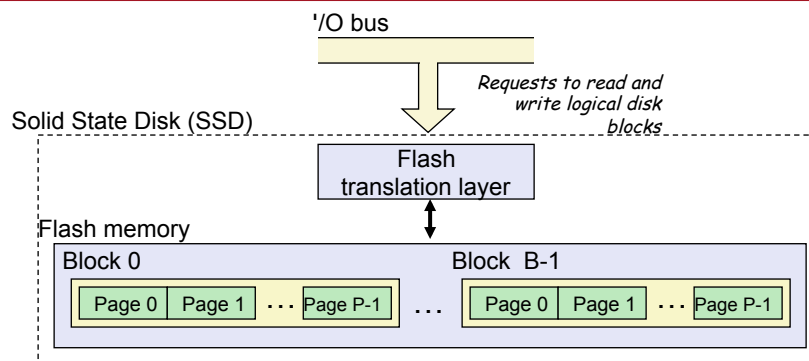
Reading a Disk Sector (2)



Reading a Disk Sector (3)



Solid State Disks (SSDs)



- Pages: 512KB to 4KB, Blocks: 32 to 128 pages
- Data read/written in units of pages.
- Page can be written only after its block has been erased
- A block wears out after 100,000 repeated writes.

SSD Performance Characteristics

Sequential read tput	250 MB/s	Sequential write tput	170 MB/s
Random read tput	140 MB/s	Random write tput	14 MB/s
Rand read access	30 us	Random write access	300 us

- Why are random writes so slow?
 - Erasing a block is slow (around 1 ms)
 - Write to a page triggers a copy of all useful pages in the block
 - Find an used block (new block) and erase it
 - Write the page into the new block
 - Copy other pages from old block to the new block

SSD Tradeoffs vs Rotating Disks

- Advantages
 - No moving parts → faster, less power, more rugged
- Disadvantages
 - Have the potential to wear out
 - Mitigated by "wear leveling logic" in flash translation layer
 - E.g. Intel X25 guarantees 1 petabyte (10¹⁵ bytes) of random writes before they wear out
 - In 2010, about 100 times more expensive per byte
- Applications
 - MP3 players, smart phones, laptops
 - Beginning to appear in desktops and servers

Storage Trends

SRAM

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	19,200	2,900	320	256	100	75	60	320
access (ns)	300	150	35	15	3	2	1.5	200

DRAM

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	8,000	880	100	30	1	0.1	0.06	130,000
access (ns)	375	200	100	70	60	50	40	9
typical size (MB)	0.064	0.256	4	16	64	2,000	8,000	125,000

Disk

Metric	1980	1985	1990	1995	2000	2005	2010	2010:1980
\$/MB	500	100	8	0.30	0.01	0.005	0.0003	1,600,000
access (ms)	87	75	28	10	8	4	3	29
typical size (MB)	1	10	160	1,000	20,000	160,000	1,500,000	1,500,000

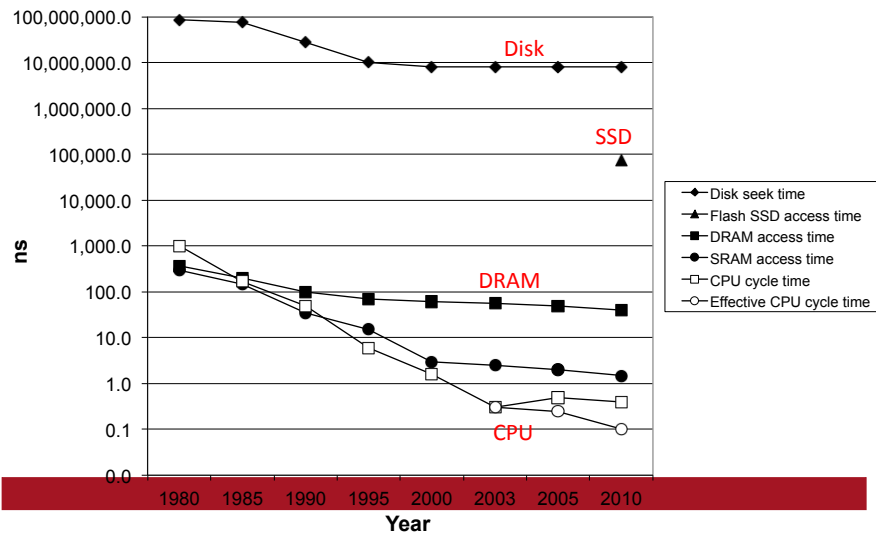
CPU Clock Rates

Inflection point in computer history
when designers hit the "Power Wall"

	1980	1990	1995	2000	2003	2005	2010	2010:1980
CPU	8080	386	Pentium	P-III	P-4	Core 2	Core i7	---
Clock rate (MHz)	1	20	150	600	3300	2000	2500	2500
Cycle time (ns)	1000	50	6	1.6	0.3	0.50	0.4	2500
Cores	1	1	1	1	1	2	4	4
Effective cycle time (ns)	1000	50	6	1.6	0.3	0.25	0.1	10,000

The CPU-Memory Gap

The gap widens between DRAM, disk, and CPU speeds.



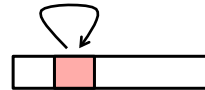
Locality

Locality

- **Principle of Locality:** Programs tend to use data and instructions with addresses near or equal to those they have used recently

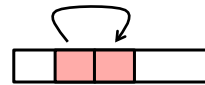
- **Temporal locality:**

- Recently referenced items are likely to be referenced again in the near future



- **Spatial locality:**

- Items with nearby addresses tend to be referenced close together in time



Locality Example

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

- Data references
 - Reference array elements in succession (stride-1 reference pattern).
 - Reference variable `sum` each iteration.
- Instruction references
 - Reference instructions in sequence.
 - Cycle through loop repeatedly.

Spatial locality

Temporal locality

Spatial locality

Temporal locality

Qualitative Estimates of Locality

- **Claim:** Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.
- **Question:** Does this function have good locality with respect to array `a`?

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}
```

Locality Example

- **Question:** Does this function have good locality with respect to array `a`?

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;

    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}
```

Locality Example

- **Question:** Can you permute the loops so that the function scans the 3-d array `a` with a stride-1 reference pattern (and thus has good spatial locality)?

```
int sum_array_3d(int a[M][N][N])
{
    int i, j, k, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < N; k++)
                sum += a[k][i][j];
    return sum;
}
```

Memory Hierarchies

- Some fundamental and enduring properties of hardware and software:
 - Fast storage technologies cost more per byte, have less capacity, and require more power (heat!).
 - The gap between CPU and main memory speed is widening.
 - Well-written programs tend to exhibit good locality.
- These fundamental properties complement each other beautifully.
- They suggest an approach for organizing memory and storage systems known as a **memory hierarchy**.

Summary

- The speed gap between *CPU*, memory and mass storage continues to widen.
 - Well-written programs exhibit a property called locality.
 - Memory hierarchies based on caching close the gap by exploiting locality.
-