



BITS Pilani
Pilani Campus

Computer Organization and Software Systems

CONTACT SESSION 2

Dr. Lucy J. Gudino
WILP & Department of CS & IS

Today's Class



Contact Hour	List of Topic Title	Text/Ref Book/external resource
3	Performance Assessment MIPS Rate Amdahl's Law	Class Slides
4	Memory Organization Storage Technologies Random Access Memory Disk Storage Solid State Disks Storage Technology Trends	T1, R2



Performance Assessment

BITS Pilani
Pilani Campus

- Kilo- (K) = 1 thousand = 10^3 and 2^{10}
- Mega- (M) = 1 million = 10^6 and 2^{20}
- Giga- (G) = 1 billion = 10^9 and 2^{30}
- Tera- (T) = 1 trillion = 10^{12} and 2^{40}
- Peta- (P) = 1 quadrillion = 10^{15} and 2^{50}
- Exa - (E) = 1 quintillion = 10^{18} and 2^{60}

Byte = a unit of storage

- 1KB = 2^{10} = 1024 Bytes \Rightarrow 2048B
- 1MB = 2^{20} = 1,048,576 Bytes
- Main memory (RAM) is measured in MB / GB
- Disk storage is measured in GB for small systems, TB for large systems.

Examples



Hertz = clock cycles per second (frequency)

- 1MHz = 1,000,000Hz
- Processor speeds are measured in MHz or GHz.

- Milli- (m) = 1 thousandth = 10^{-3}
- Micro- (μ) = 1 millionth = 10^{-6}
- Nano- (n) = 1 billionth = 10^{-9}
- Pico- (p) = 1 trillionth = 10^{-12}
- Femto- (f) = 1 quadrillionth = 10^{-15}

Examples

- Millisecond = 1 thousandth of a second
 - Hard disk drive access times are often 10 to 20 milliseconds.
- Nanosecond = 1 billionth of a second
 - Main memory access times are often 50 to 70 nanoseconds.
- Micron (micrometer) = 1 millionth of a meter
 - Circuits on computer chips are measured in microns.

Important Terms



- **Execution time** : The total time required for the computer to complete a task, including disk accesses, memory accesses, I/O activities, operating system overhead, CPU execution
- **Throughput or bandwidth** : number of tasks completed per unit time.

Example



What changes to a computer system will increase throughput, decrease execution time, or both?

1. Replacing the processor in a computer with a faster version
2. Adding additional processors of same type to a system,
that is, it uses multiple processors for separate tasks

Contd...



- Relationship between Performance and execution time of Computer X

$$\text{Performance}_x = \frac{1}{\text{Execution time}_x}$$

- if the performance of X is greater than the performance of Y, we have

$$\text{Performance}_x > \text{Performance}_y$$

$$\frac{1}{\text{Execution time}_x} > \frac{1}{\text{Execution time}_y}$$

$$\text{Execution time}_y > \text{Execution time}_x$$

Contd...



- Quantitative performance analysis
 - Computer X is "n" times faster than Computer Y

$$\frac{\text{Performance}_x}{\text{Performance}_y} = n$$

$$\frac{1}{E_x} \quad \frac{E_y}{E_x}$$

$$\left[\frac{\text{Performance}_x}{\text{Performance}_y} = \frac{\text{Execution time}_y}{\text{Execution time}_x} = n \right]$$

- If performance of X is n times better than Y, then the execution time on Y is n times longer than it is on X

Example



- If computer A runs a program in 10 seconds and computer B runs the same program in 15 seconds, how much faster is A than B?

$$\frac{\text{Performance}_A}{\text{Performance}_B} = \frac{\text{Execution time}_B}{\text{Execution time}_A} = n$$
$$= \frac{15}{10} \Rightarrow 1.5$$

- Computer A is therefore 1.5 times faster than B.

CPU performance and its factors



$$\frac{\text{Performance}_x}{\text{Performance}_y} = \frac{\text{Execution time}_y}{\text{Execution time}_x} = n$$

- CPU execution time for a program:

$$\text{CPU execution time for a program} = \text{CPU clock cycles for a program} \times \text{Clock cycle time}$$



$$\text{CPU execution time for a program} = \frac{\text{CPU clock cycles for a program}}{\text{Clock rate}}$$

Example



- Our favorite program runs in 10 seconds on computer A, which has a 2 GHz clock. We are trying to help a computer designer build a computer, B, which will run this program in 6 seconds. The designer has determined that a substantial increase in the clock rate is possible, but this increase will affect the rest of the CPU design, causing computer B to require 1.2 times as many clock cycles as computer A for this program. What clock rate should we tell the designer to target?

$$\text{CPU execution time for a program} = \frac{\text{CPU clock cycles for a program}}{\text{Clock rate}}$$

Computer A

Execution Time_A = 10s

Clock Rate_A = 2×10^9 Hz

CPU Clock Cycle_A = ?

$$10s = \frac{CC_A}{2 \times 10^9}$$

$$20 \times 10^9 = CC_A$$

Computer B

Execution Time_B = 6s

CPU Clock Cycles_B = $1.2 \times$ Clock Cycle_A

Clock Rate B = ?

$$6 = \frac{1.2 \times CC_A}{CR_B}$$

$$\begin{aligned} CR_B &= \frac{1.2 \times 20 \times 10^9}{6} \\ &= 0.2 \times 20 \times 10^9 \\ &= 4.0 \times 10^9 \end{aligned}$$

Let's first find the number of clock cycles required for the program on A:

$$\text{CPU time}_A = \frac{\text{CPU clock cycles}_A}{\text{Clock rate}_A}$$

$$10 \text{ seconds} = \frac{\text{CPU clock cycles}_A}{2 \times 10^9 \frac{\text{cycles}}{\text{second}}}$$

$$\text{CPU clock cycles}_A = 10 \text{ seconds} \times 2 \times 10^9 \frac{\text{cycles}}{\text{second}} = 20 \times 10^9 \text{ cycles}$$

CPU time for B can be found using this equation:

$$\text{CPU time}_B = \frac{1.2 \times \text{CPU clock cycles}_A}{\text{Clock rate}_B}$$

$$6 \text{ seconds} = \frac{1.2 \times 20 \times 10^9 \text{ cycles}}{\text{Clock rate}_B}$$

$$\text{Clock rate}_B = \frac{1.2 \times 20 \times 10^9 \text{ cycles}}{6 \text{ seconds}} = \frac{0.2 \times 20 \times 10^9 \text{ cycles}}{\text{second}} = \frac{4 \times 10^9 \text{ cycles}}{\text{second}} = 4 \text{ GHz}$$

Instruction Performance



- CPI: Clock cycles Per Instruction
 - Average number of clock cycles per instruction for a program or program fragment.

$$\text{CPU clock cycles} = \text{Instructions for a program} \times \text{Average clock cycles per instruction}$$

Example



Computer A has a clock cycle time of 250 ps and a CPI of 2.0 for some program, and computer B has a clock cycle time of 500 ps and a CPI of 1.2 for the same program. Which computer is faster for this program and by how much?

Solution

Computer A has a clock cycle time of 250 ps and a CPI of 2.0 for some program, and computer B has a clock cycle time of 500 ps and a CPI of 1.2 for the same program. Which computer is faster for this program and by how much?



- The number of processor clock cycles for each computer

$$\text{CPU clock cycles}_A = I \times 2.0$$

$$\text{CPU clock cycles}_B = I \times 1.2$$

- Execution time for each computer

$$\text{Execution time} = \text{CPU clock cycles} \times \text{Clock cycle time}$$

$$\text{Execution time}_A = I \times 2.0 \times 250 \text{ ps} = 500 \times I \text{ ps}$$

$$\text{Execution time}_B = I \times 1.2 \times 500 \text{ ps} = 600 \times I \text{ ps}$$

- Comparison:

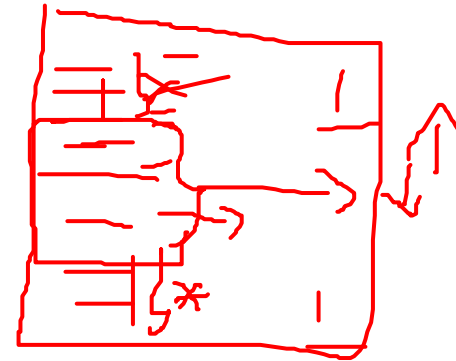
$$\frac{\text{CPU performance}_A}{\text{CPU performance}_B} = \frac{\text{Execution time}_B}{\text{Execution time}_A} = \frac{600 I \text{ ps}}{500 I \text{ ps}} = 1.2$$

Amdahl's Law



- proposed by Gene Amdahl in 1967
- deals with the potential speedup of a program using multiple processors compared to a single processor

$$\text{Speedup} = \frac{\text{Performance after enhancement}}{\text{Performance before enhancement}} = \frac{\text{Execution time before enhancement}}{\text{Execution time after enhancement}}$$



Amdahl's Law



$$\text{Speedup} = \frac{\text{Performance after enhancement}}{\text{Performance before enhancement}} = \frac{\text{Execution time before enhancement}}{\text{Execution time after enhancement}}$$

$$S = \frac{1}{(1-f) + \frac{f}{k}}$$

Handwritten red annotations: a checkmark to the left of the equation, a red 'f' with a checkmark below the denominator, and a red arrow pointing from 'f' to 'P'.

S=Speedup,
f=fraction of time enhancement,
k=speedup of the faster component

$$\underline{S} = \frac{1}{(1-P) + \frac{P}{k}}$$

Handwritten red annotations: 'no of times' written above the fraction $\frac{P}{k}$, and ' $\frac{1}{6}$ ' written below it.

Amdahl's Law

$$S = \frac{1}{(1-P) + \frac{P}{K}}$$



If 90% of a program is speeded up to run 10 times faster $f=0.9$ and $k=10$

Overall speedup is $1/(1-0.9) + (0.9/10) =$
 $1/(0.1 + 0.09) = 1/(0.19) = \underline{\underline{5.26}}$

Making 80% of a program run 20% faster

$f=0.80$ and $k=1.2$ $K = 100 + \frac{20}{0.8} \Rightarrow 1 + \frac{20}{100}$
 $1/(1-0.8) + (0.8/1.2) =$
 $1/(0.2 + 0.8/1.2) = 1/(0.2 + 0.66) = 1/0.866 = \underline{\underline{1.154}} \Rightarrow 1.2 //$

Example



On a large system CPU upgrade makes it faster by 50% for INR 10,000. A disk drive upgrade of INR 7000 speeds it up by 150%. Evaluate the speedups? Processes spend 70% in CPU and 30% waiting Disk drives.

Processor upgrade

Disk Drive upgrade

$$f = 0.70, \quad k = 1.5, \quad S = \frac{1}{(1 - 0.7) + 0.7/1.5} = 1.304$$

$$f = 0.30, \quad k = 2.5, \quad S = \frac{1}{(1 - 0.3) + 0.3/2.5} = 1.219$$

30% improvement

22% Improvement

CPU-30 % improvement -faster by 50%
---so 1% increment is INR 10000/30=INR 333

Handwritten calculation: $1.5 \rightarrow 1.5 + 1 \rightarrow 2.5$

DISK DRIVE- 22% improvement – speeds up 150%---so a 1% increment is INR 7000/22=INR=318

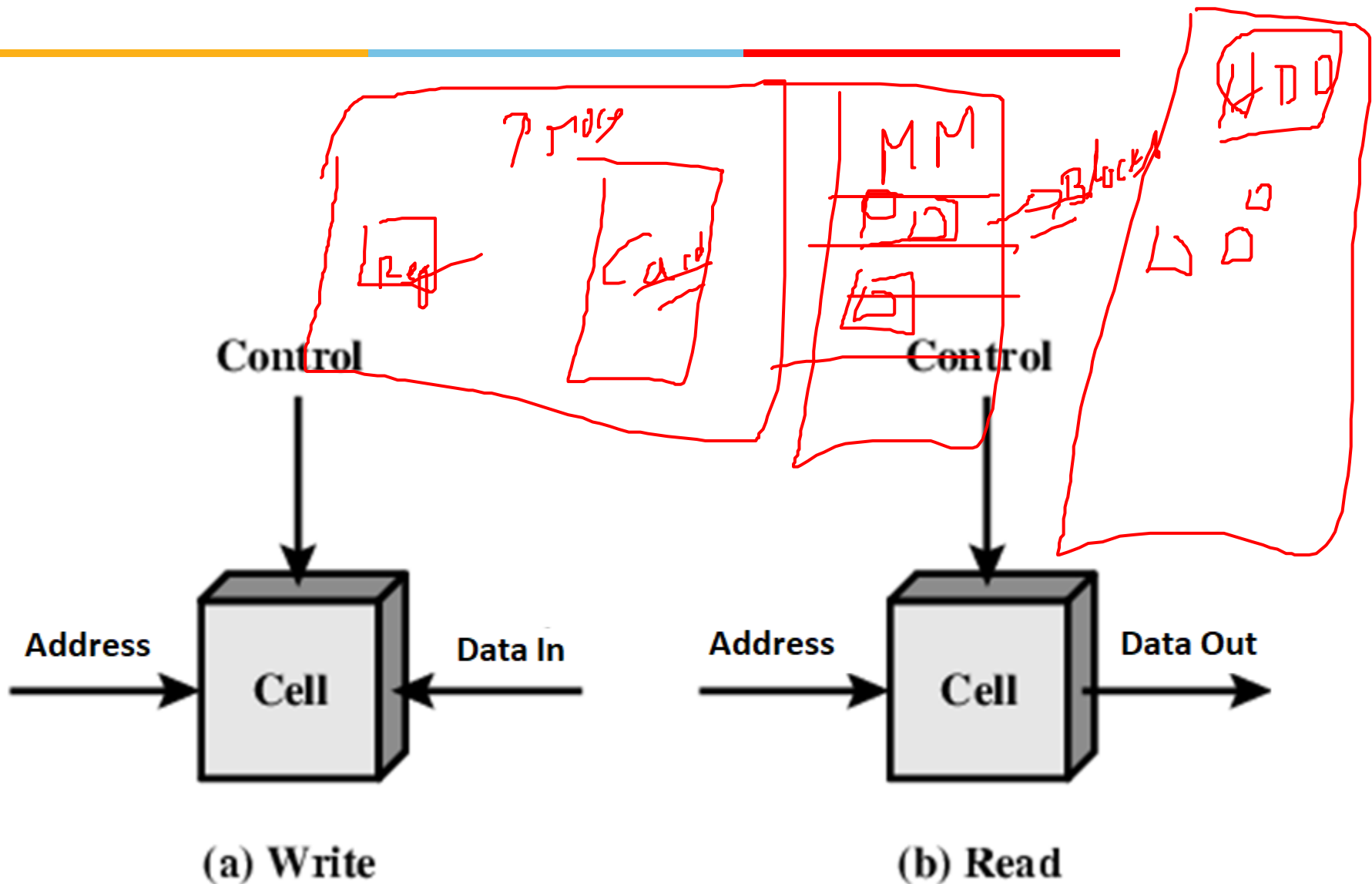
Each 1% of improvement for the processor costs INR333, and for the disk a 1% improvement costs INR318. "Is cost/performance the most important metric?"



BITS Pilani
Pilani Campus

Memory Organization

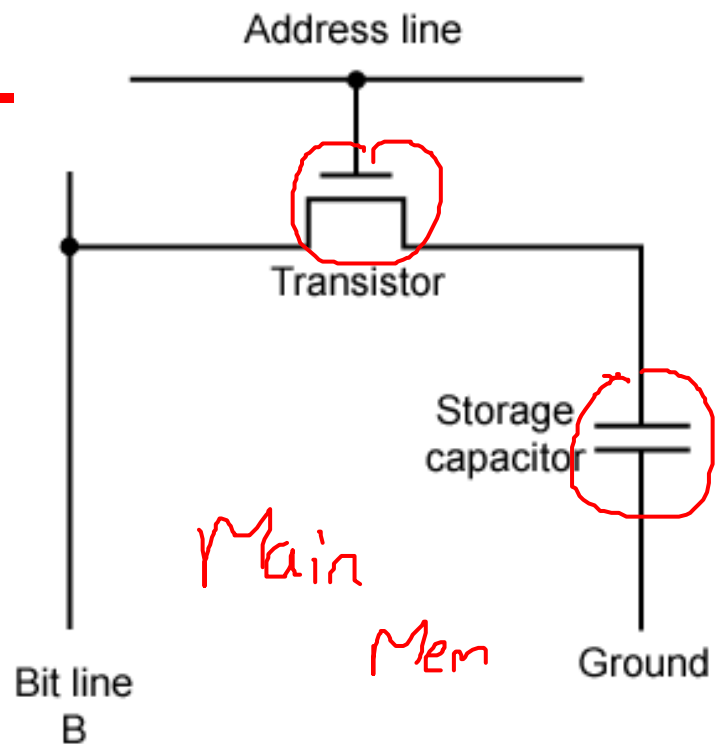
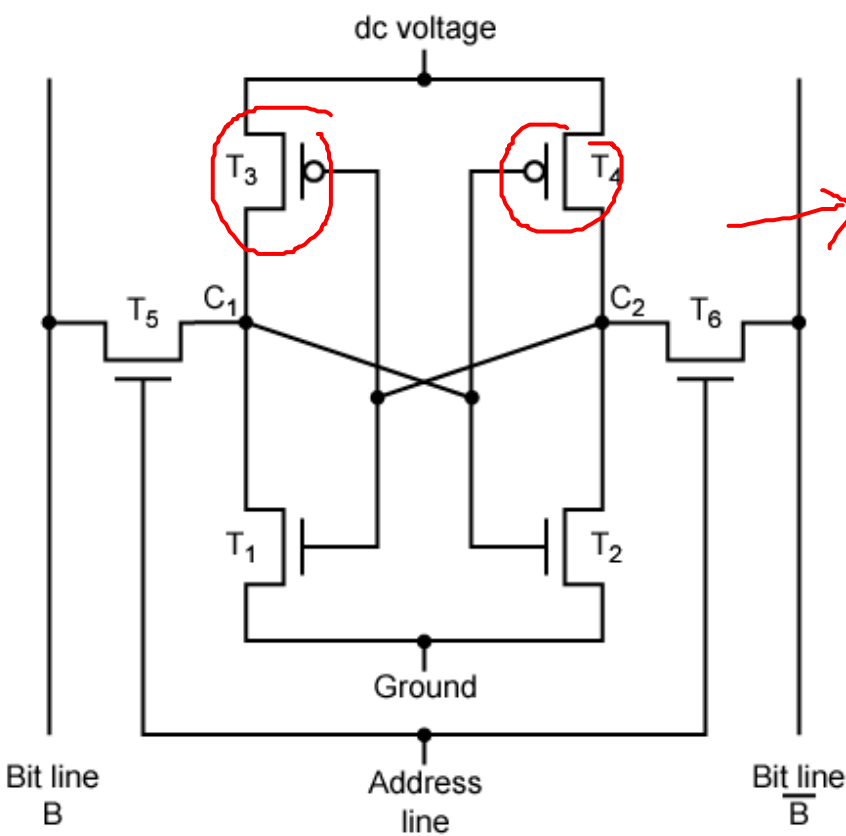
Semiconductor Memory



Random-Access Memory (RAM)

- Key features
 - **RAM** is traditionally packaged as a chip.
 - Basic storage unit is normally a **cell** (one bit per cell).
 - Multiple RAM chips form a memory.
- RAM comes in two varieties:
 - ☒ SRAM (Static RAM)
 - ☒ DRAM (Dynamic RAM)
- SRAM and DRAM are volatile memories
 - Lose information if powered off.

SRAM vs DRAM Summary



	Trans. per bit	Access time	Needs refresh?	Needs EDC?	Cost	Applications
SRAM	4 to 6	1X	No	Maybe	100x	Cache
DRAM	1	10X	Yes	Yes	1X	Main memories, frame buffers

Read Only Memory



- Permanent Storage and Nonvolatile Memories
- Read Only Memory Variants:
 - 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 (block-level) erase capability
 - Wears out after about 100,000 erasing
- Firmware

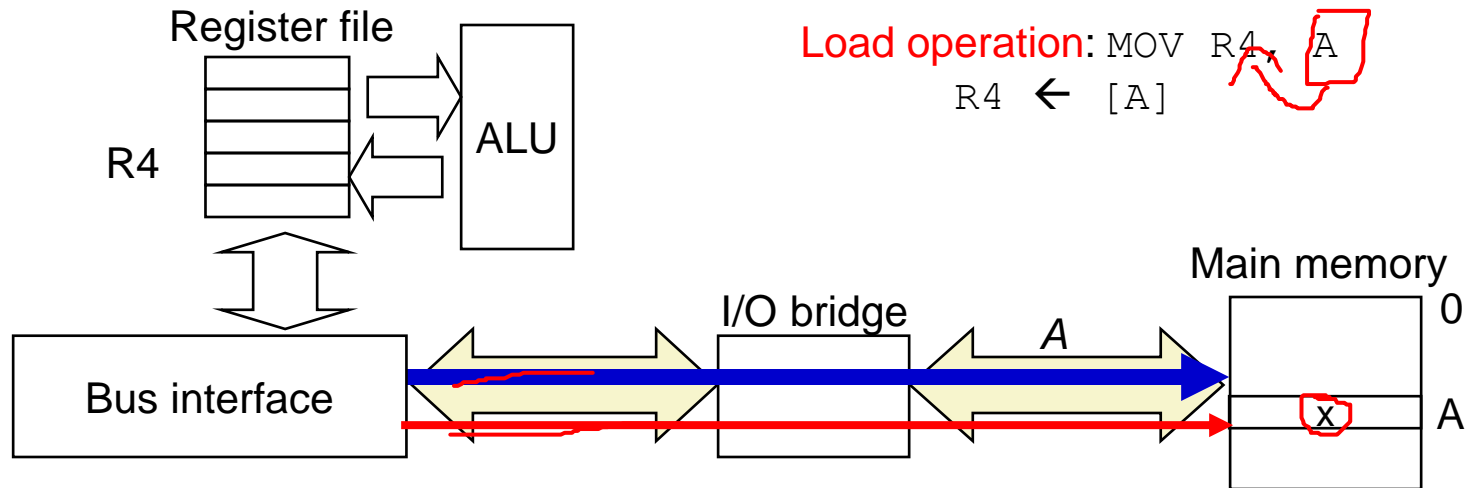
Applications

- Storing fonts for printers
- Storing sound data in musical instruments
- Video game consoles
- Implantable Medical devices.
- High definition Multimedia Interfaces(HDMI)
- BIOS chip in computer
- Program storage chip in modem, video card and many electronic gadgets, controllers for disks, network cards,

Memory Read Operation (1)



CPU places **address A** and then **read control signal** on the memory bus

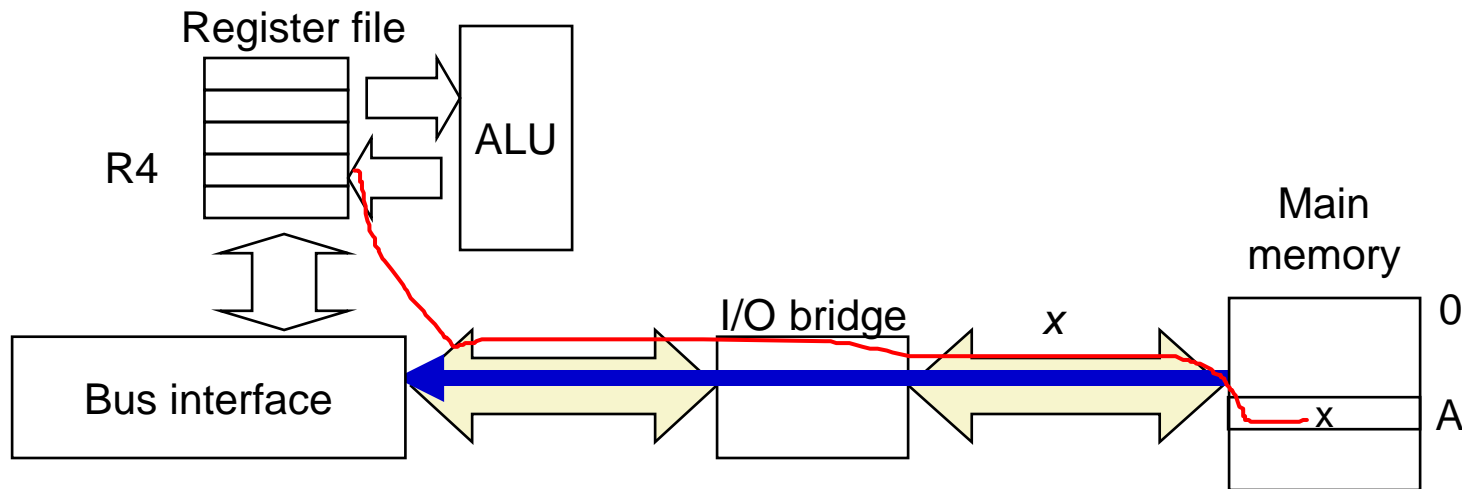


Memory Read Operation (2)



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

Load operation: `MOV R4, A`
 $R4 \leftarrow [A]$

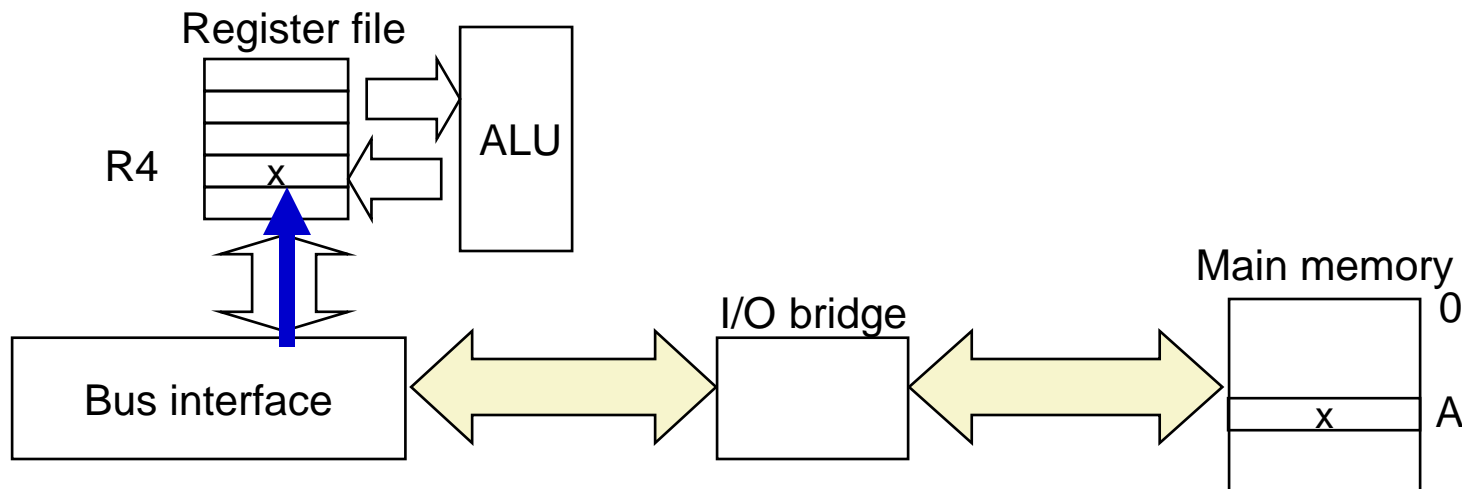


Memory Read Operation (3)



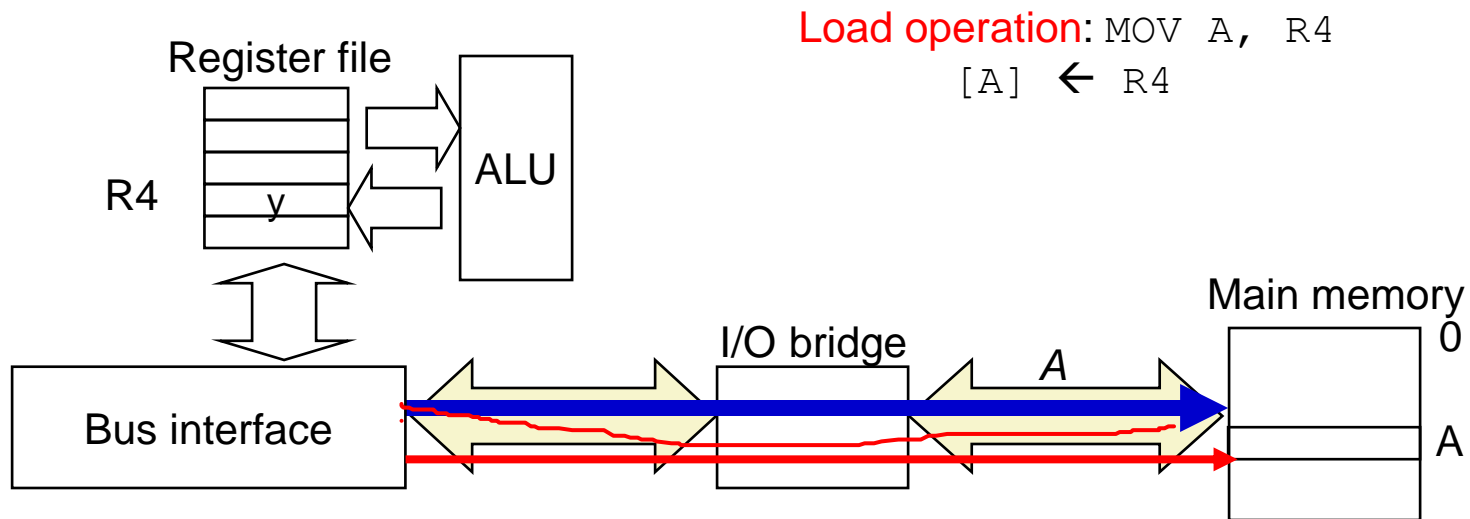
CPU read word x from the bus and copies it into register R4.

Load operation: `MOV R4, A`
 $R4 \leftarrow [A]$



Memory Write Operation (1)

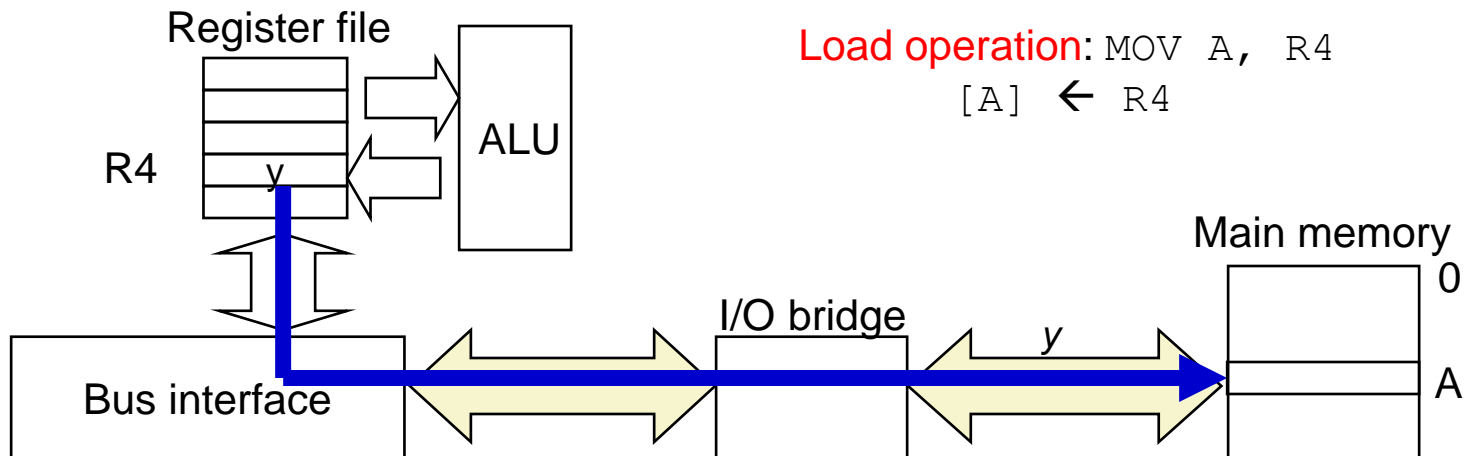
CPU places **address A** and **WRITE** control signal on bus.
Main memory reads them and waits for the corresponding data word to arrive.



Memory Write Operation (2)

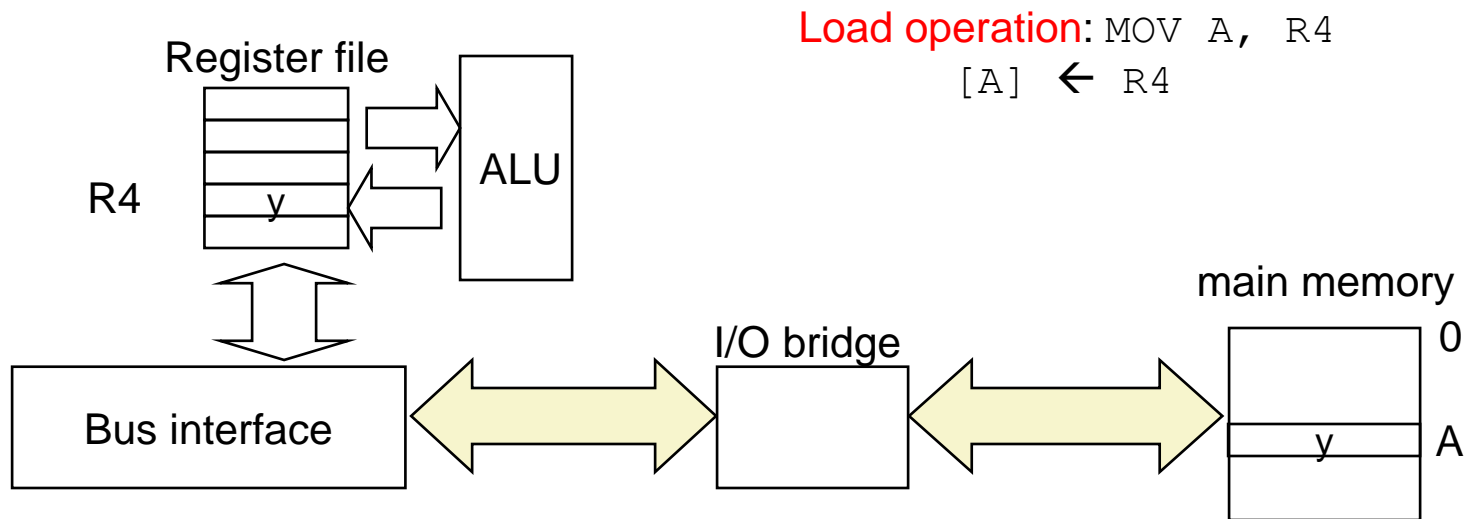


CPU places data word y on the bus



Memory Write Operation (3)

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



Magnetic Disk Drive

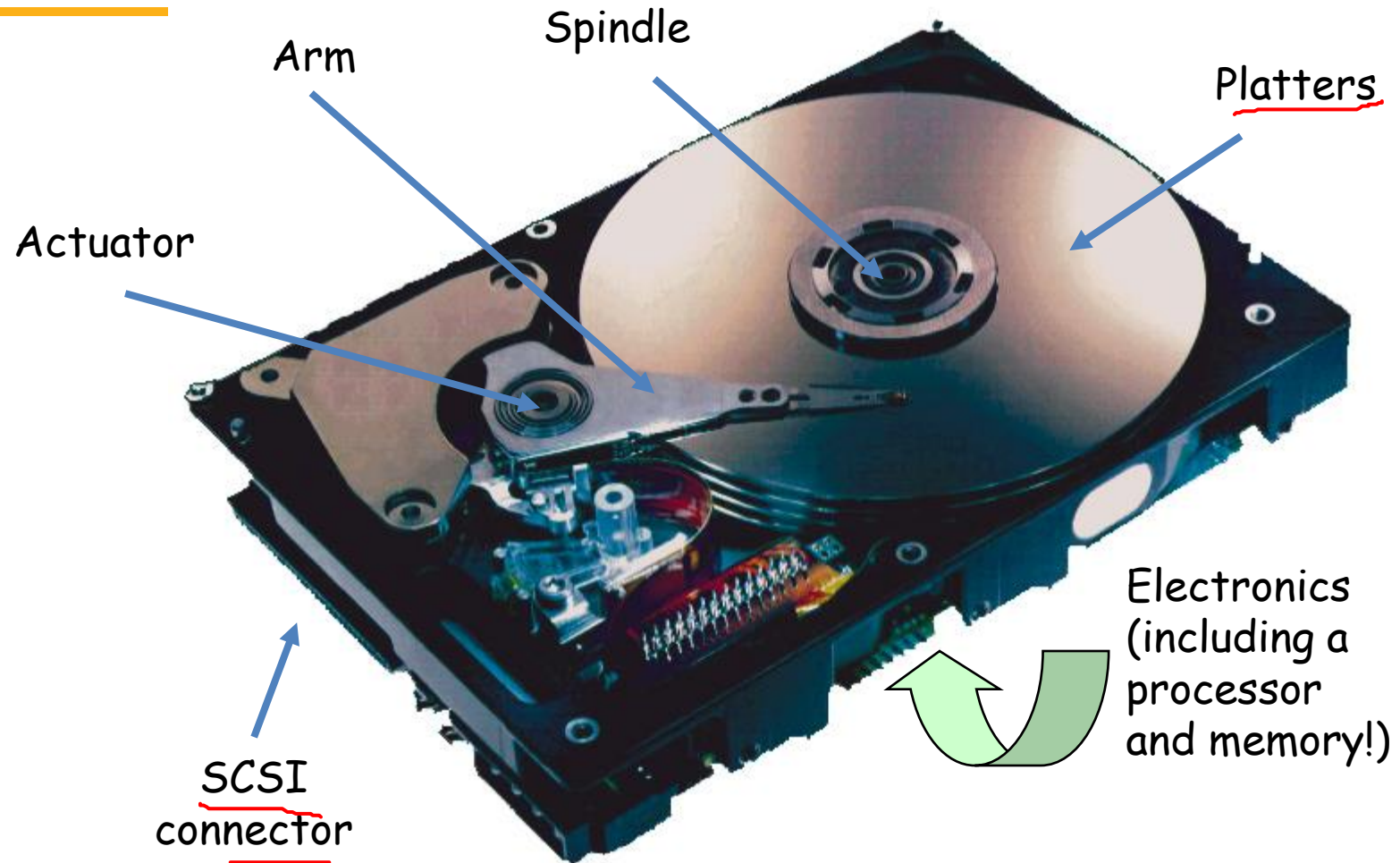
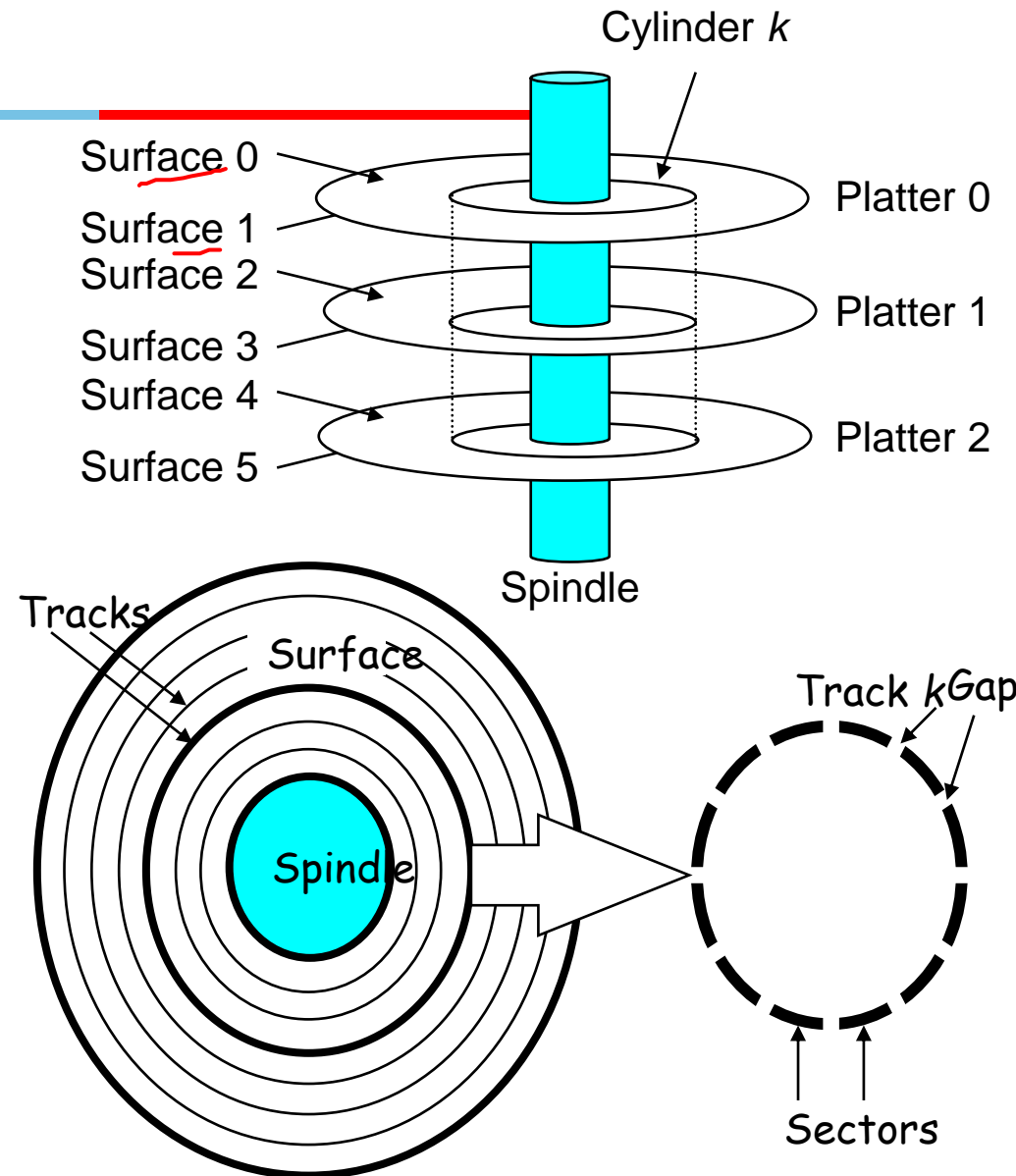


Image courtesy of Seagate Technology

Disk Geometry



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

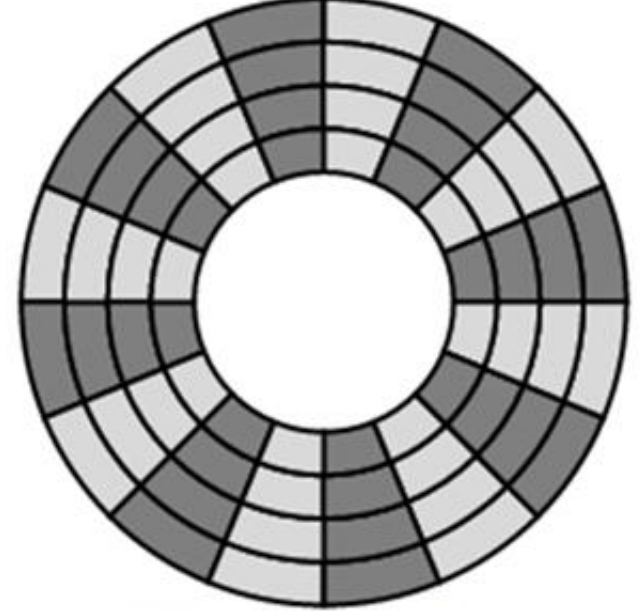


Disk Capacity

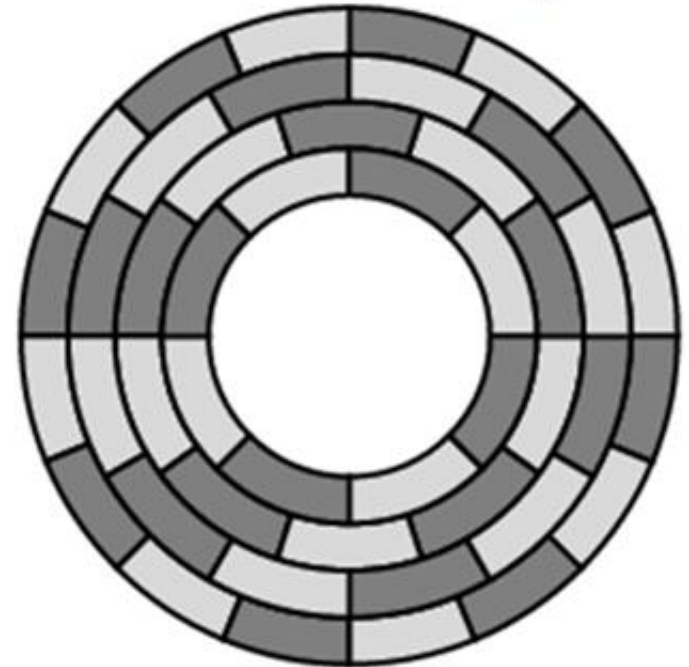
- **Capacity**: maximum number of bits that can be stored.
 - Vendors express capacity in units of gigabytes (GB /TB), where $1 \text{ GB} = 2^{30} \text{ Bytes}$, $1 \text{ TB} = 2^{40} \text{ Bytes}$,
- Capacity is determined by these technology factors:
 - **Recording density** (bits/in): number of 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.

Recording zones

- 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, outer zones have more sectors/track than inner zones.
 - So we use **average** number of sectors/track when computing capacity.



Without Recording Zones



With Recording Zones

Computing Disk Capacity



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

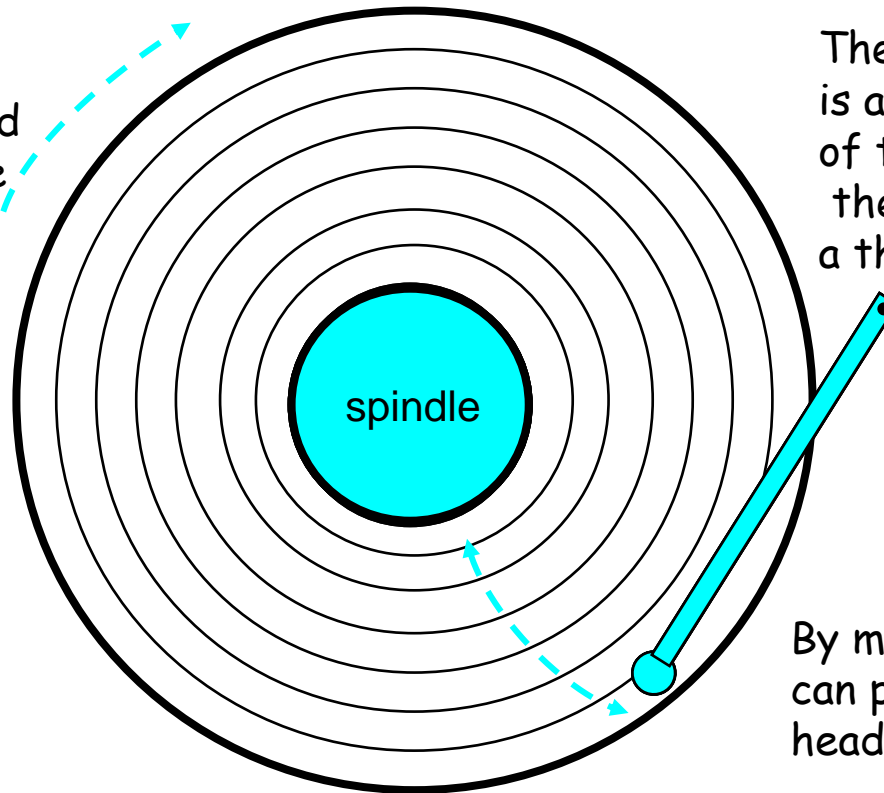
- Example:

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

- Capacity = $512 \times 300 \times 20000 \times 2 \times 5$
= $30,720,000,000$
= 28.61 GB
 $1024 \times 1024 \times 1024$

Disk Operation (Single-Platter View)

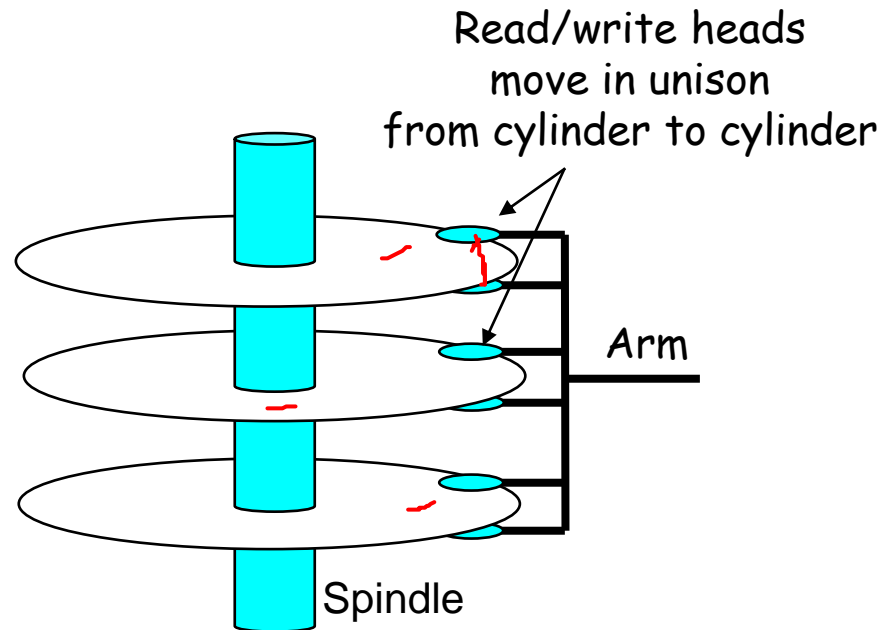
The disk surface spins at a fixed rotational rate.



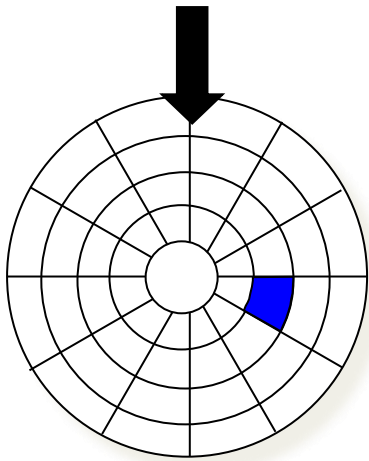
The read/write head is attached to the end of the arm and flies over the disk surface on a thin cushion of air.

By moving radially, the arm can position the read/write head over any track.

Disk Operation (Multi-Platter View)

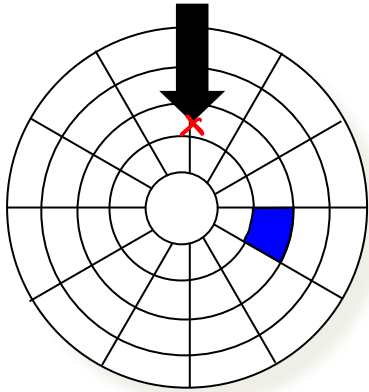


Disk Access



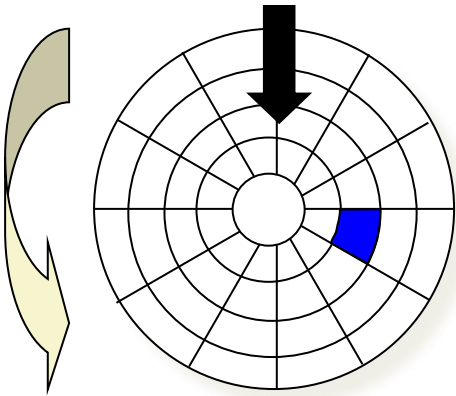
Need to access a sector colored in blue

Disk Access



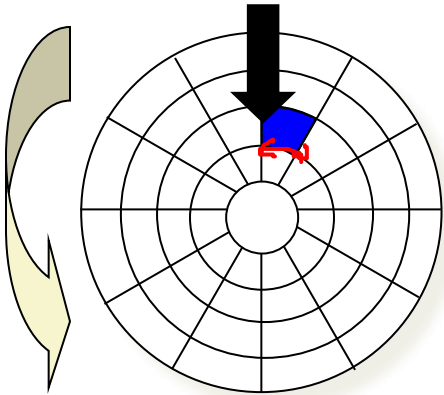
Head in position above a track

Disk Access



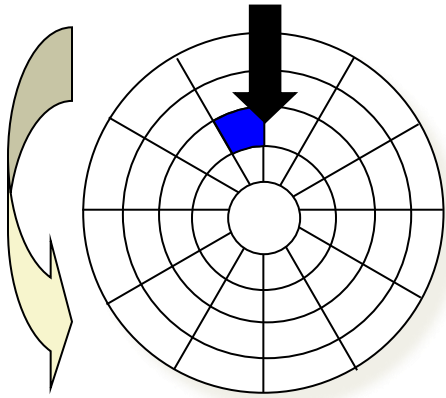
Rotate the platter in counter-clockwise direction

Disk Access - Read



About to read blue sector

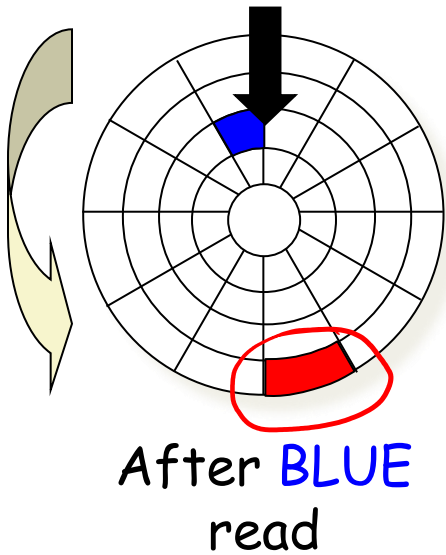
Disk Access - Read



After BLUE
read

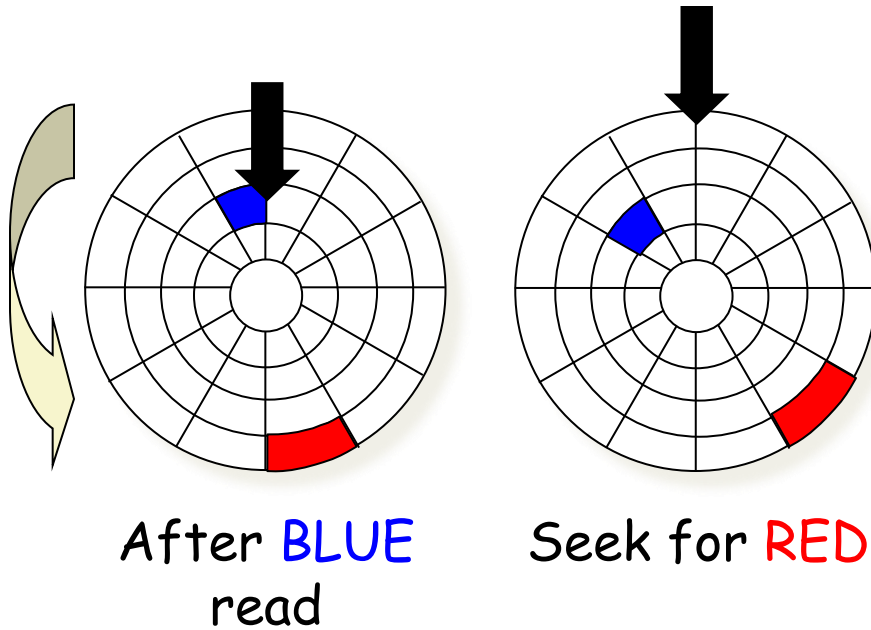
After reading blue sector

Disk Access - Read



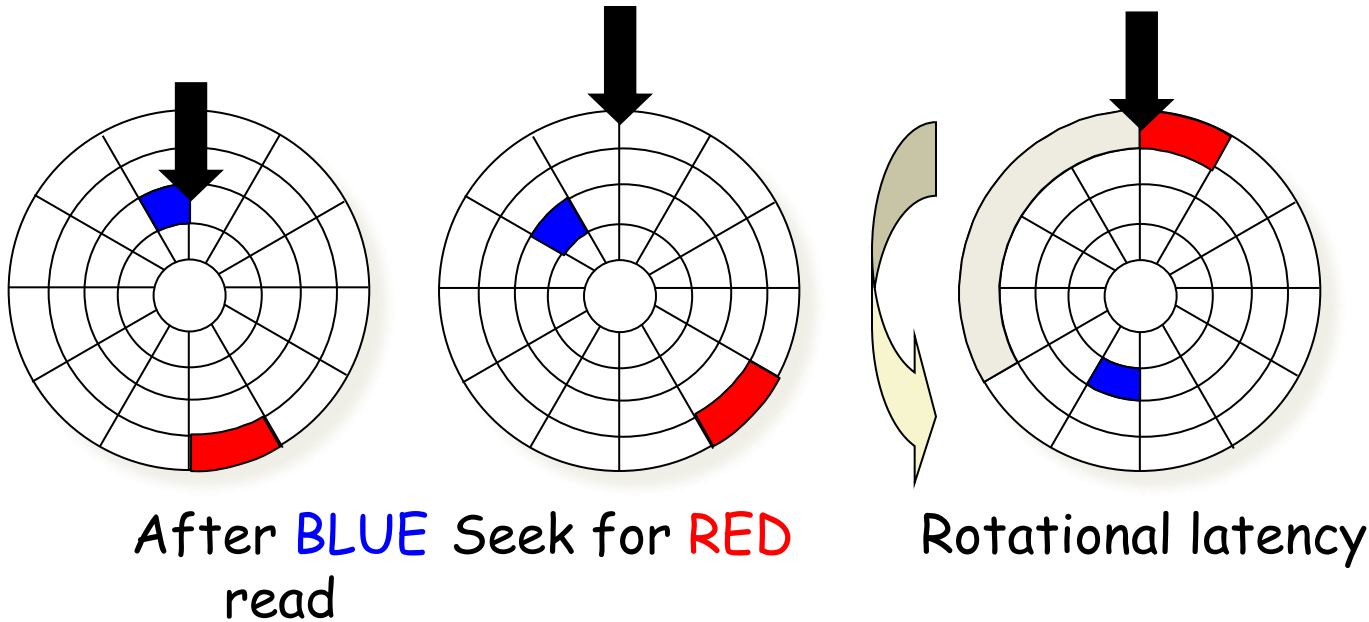
Red request scheduled next

Disk Access - Seek



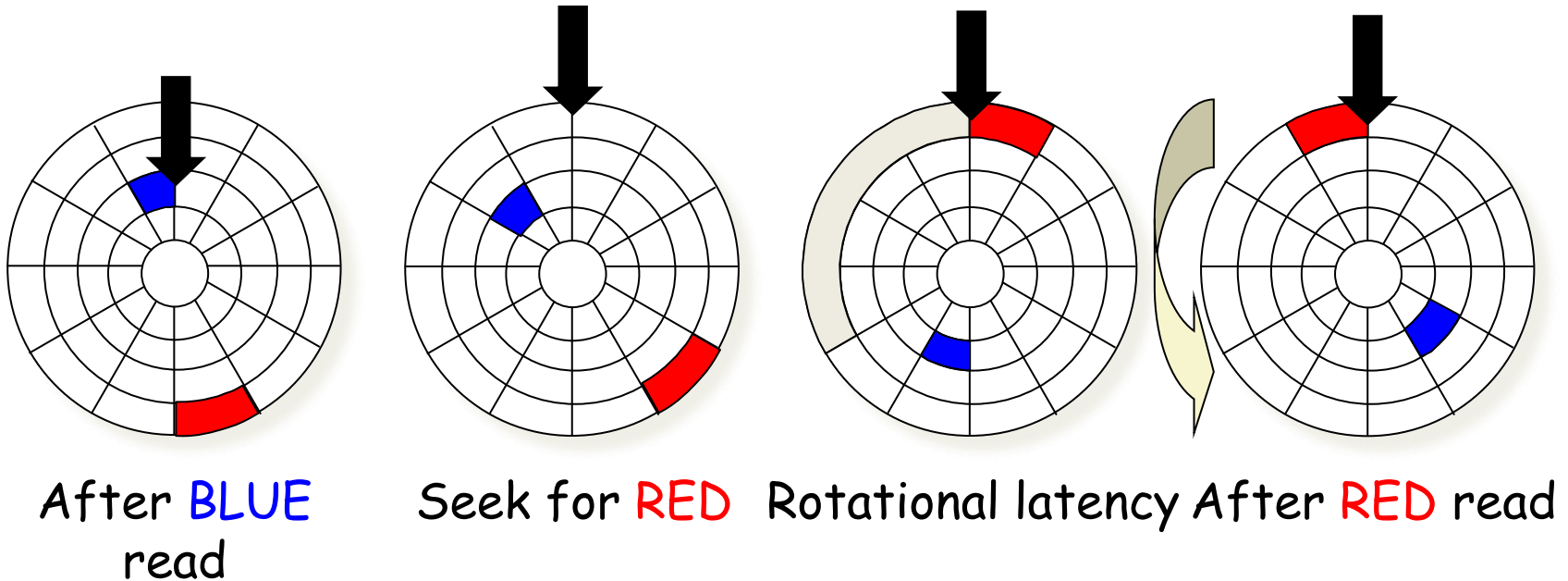
Seek to red's track

Disk Access - Rotational Latency



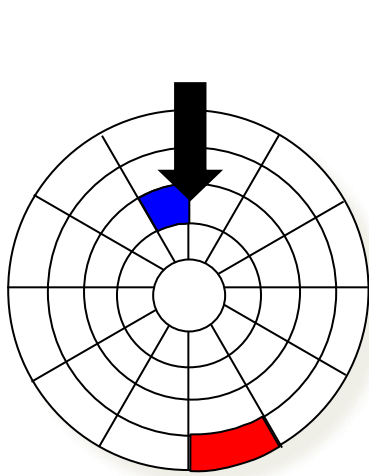
Wait for red sector to rotate around

Disk Access - Read



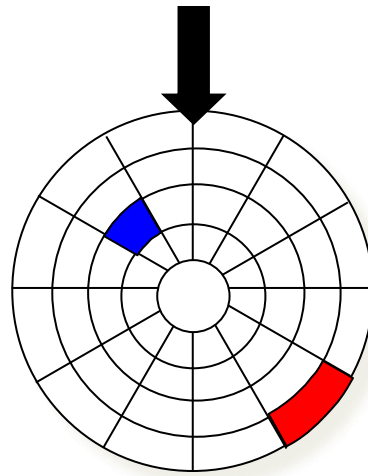
Complete read of red

Disk Access - Access Time Components



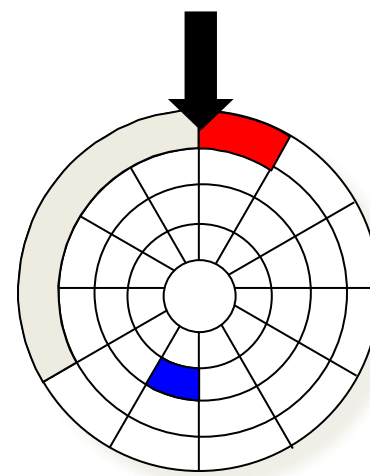
After **BLUE**
read

↑
Data transfer



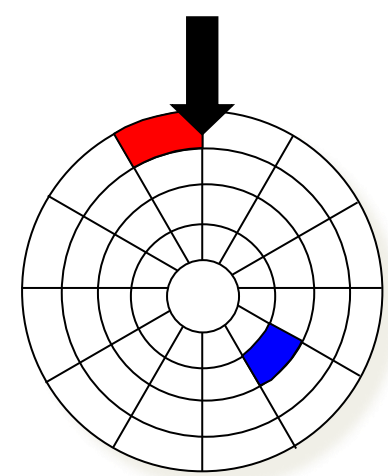
Seek for **RED**

↑
Seek



Rotational latency After **RED** read

↑
Rotational
latency



After **RED** read

↑
Data transfer

S R S-

1 2 3 4

Disk Access Time

- Average time to access some target sector given by :
 - $T_{\text{access}} = T_{\text{avg seek}} + T_{\text{avg rotation}} + T_{\text{avg transfer}}$
- **Seek time** ($T_{\text{avg seek}}$)
 - Time to position heads over cylinder containing target sector.
 - Typical $T_{\text{avg seek}}$ is 3–9 ms
- **Rotational latency** ($T_{\text{avg rotation}}$)
 - Time waiting for first bit of target sector to pass under r/w head.
 - $T_{\text{avg rotation}} = \frac{1}{2r}$, where r is rotation Speed in revolution per Second
 - Typical $T_{\text{avg rotation}} = 7200 \text{ RPMs} = 7200/60 \text{ RPS}$
- **Transfer time** ($T_{\text{avg transfer}}$)
 - Time to read the bits in the target sector.
 - $T_{\text{avg transfer}} = b/rN$, where b is the number of bytes to be transferred and N is the average number of bytes on a track

Disk Access Time Example

Given:

- Rotational rate = 7,200 RPM
- Average seek time = 9 ms.
- Avg # sectors/track = 400.
- 512 bytes per sector

$$\frac{1}{2 \times 7200} \Rightarrow \frac{60}{2 \times 7200}$$

Derived:

- Tavg rotation = $1/2r$ = $1/2 \times (60 \text{ secs}/7200 \text{ RPM})$
= $0.00416 = 4.16\text{ms.}$
- Tavg transfer = b/rN
= $512 \times 60/7200 \times 1/(400 \times 512)$
= 0.02 ms
- Taccess = $9 \text{ ms} + 4.16\text{ms} + 0.02 \text{ ms} = 13.18\text{ms}$

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.