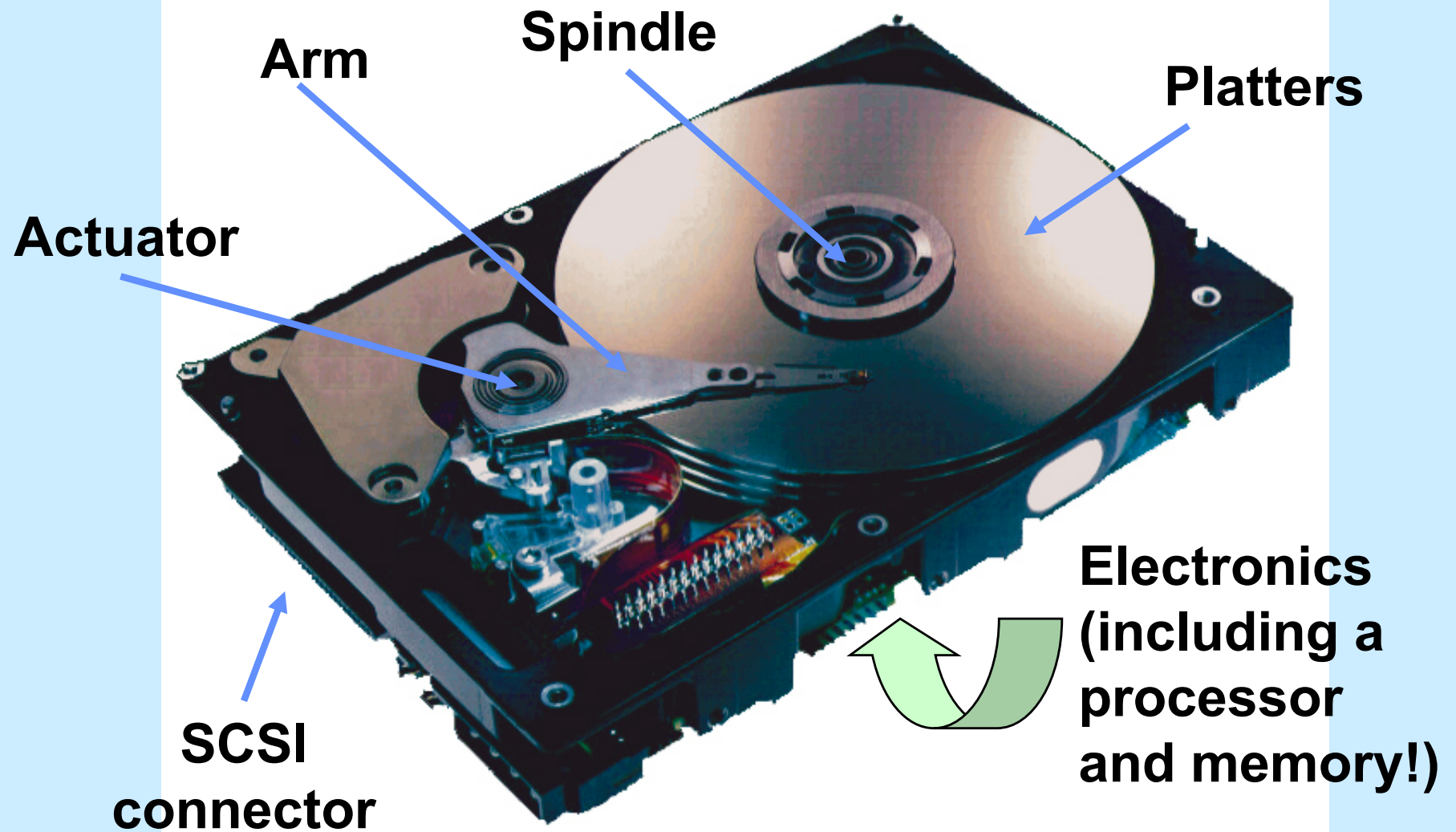


# CS 33

## Memory Hierarchy II

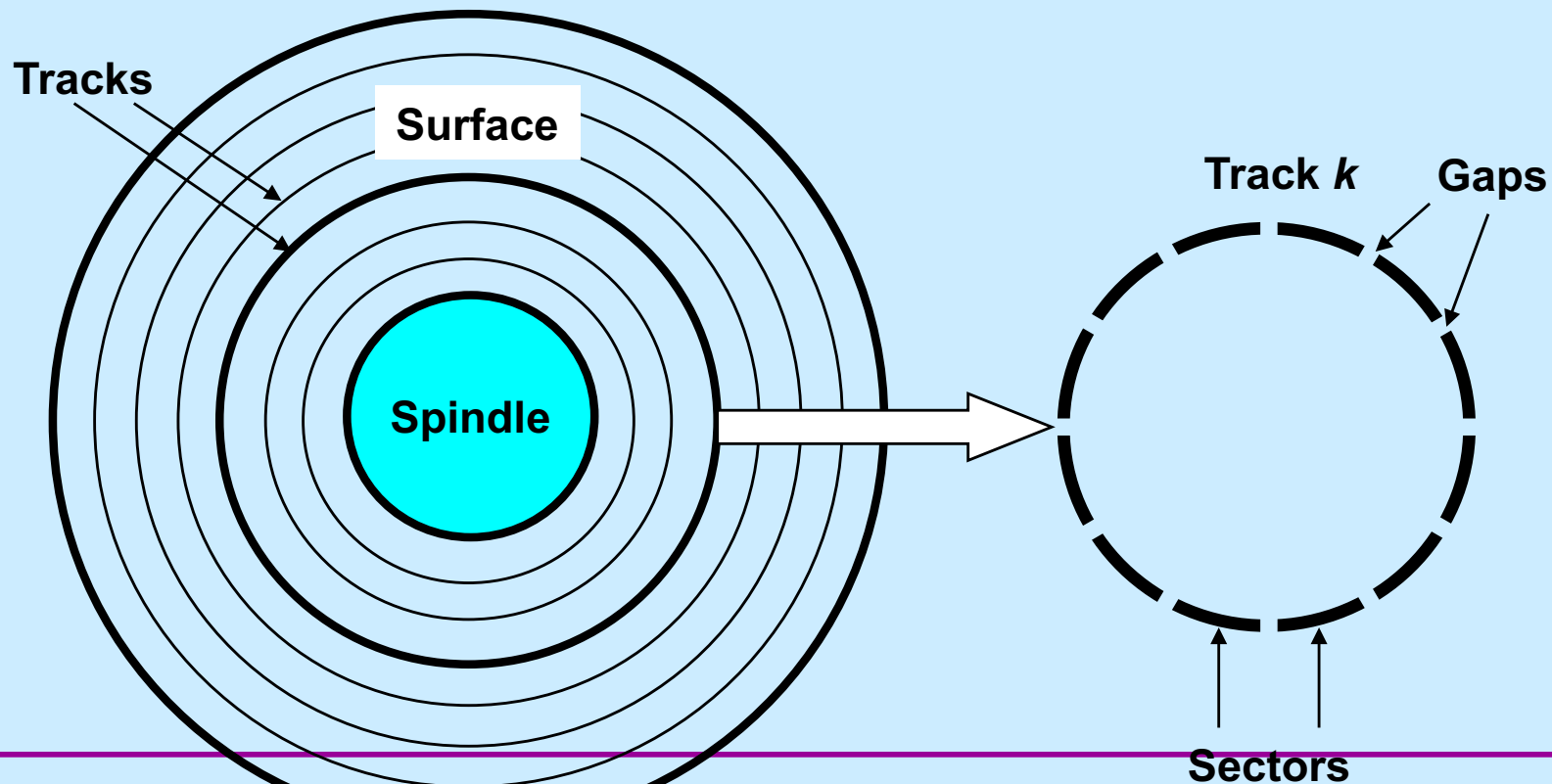
# What's Inside A Disk Drive?



*Image courtesy of Seagate Technology*

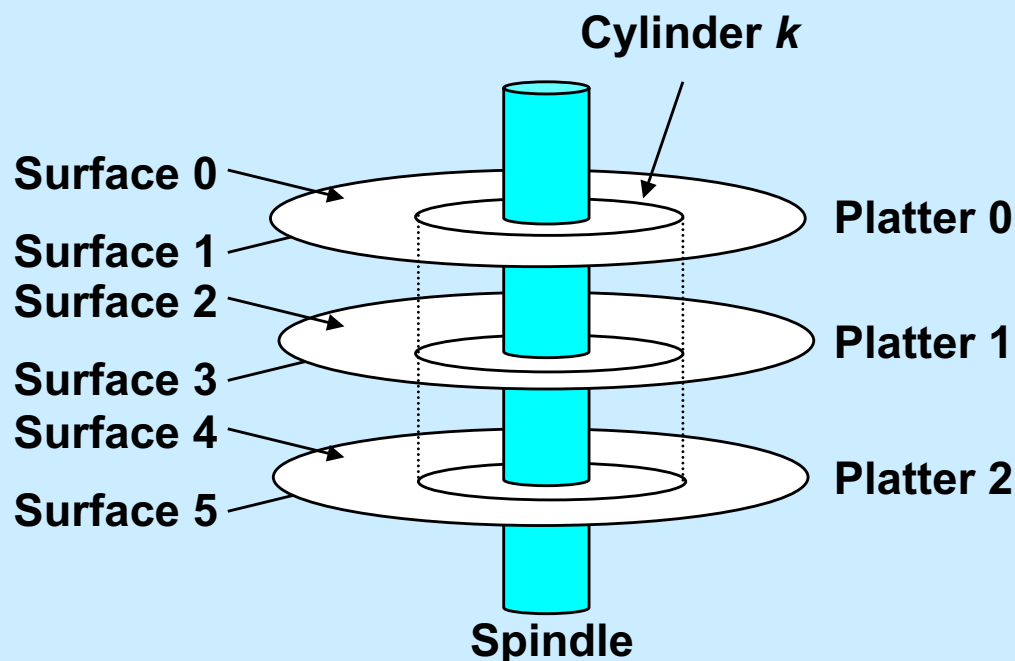
# Disk Geometry

- Disks consist of **platters**, each with two **surfaces**
- Each surface consists of concentric rings called **tracks**
- 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
  - capacity expressed in units of gigabytes (GB), where  $1 \text{ GB} = 2^{30} \text{ Bytes} \approx 10^9 \text{ Bytes}$
- 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/in<sup>2</sup>): 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

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

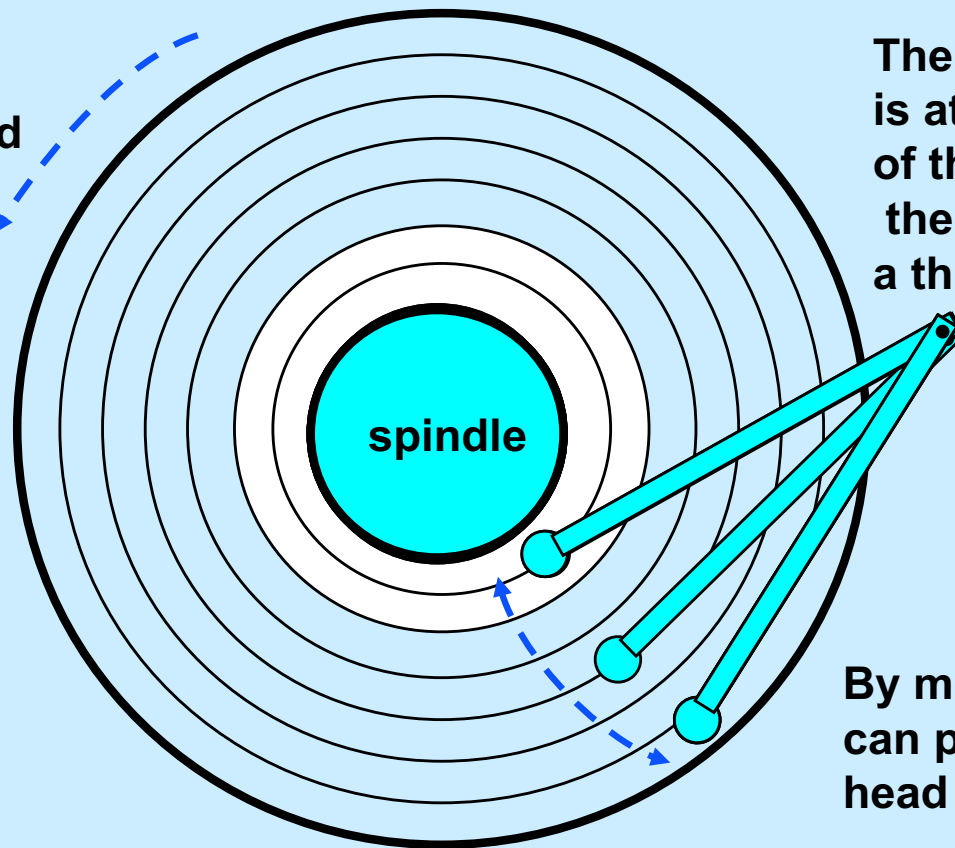
**Example:**

- 512 bytes/sector
- 600 sectors/track (on average)
- 40,000 tracks/surface
- 2 surfaces/platter
- 5 platters/disk

$$\begin{aligned}\text{Capacity} &= 512 \times 600 \times 40000 \times 2 \times 5 \\ &= 122,880,000,000 \\ &= 113.88 \text{ GB}\end{aligned}$$

# Disk Operation (Single-Platter View)

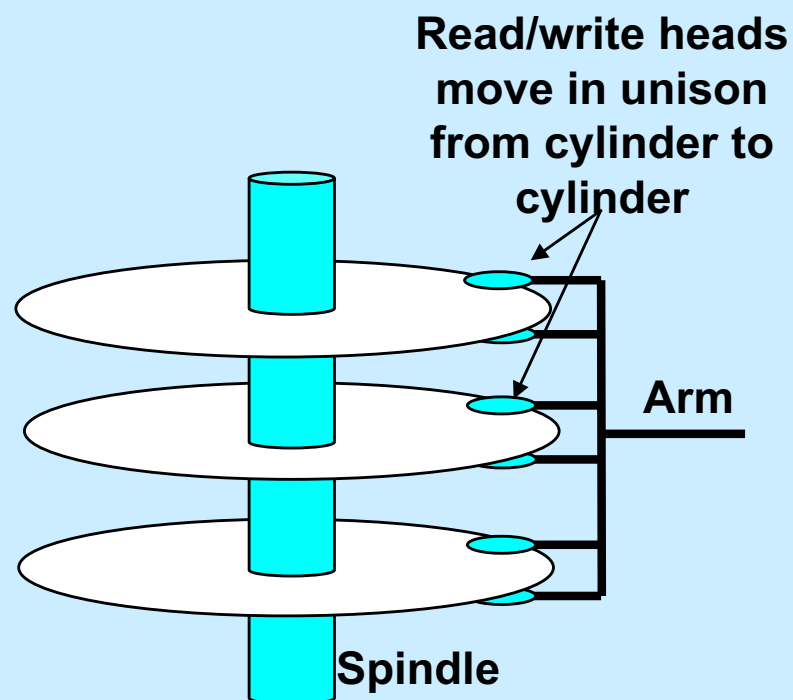
The disk surface spins at a fixed rotational rate



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

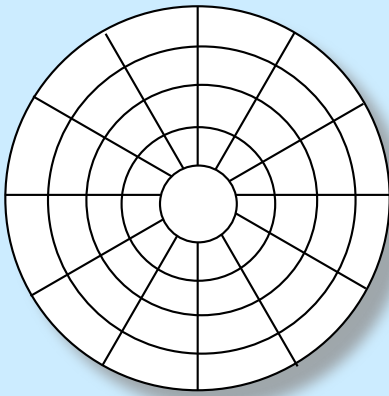
By moving radially, the arm can position the read/write head over any track

# Disk Operation (Multi-Platter View)





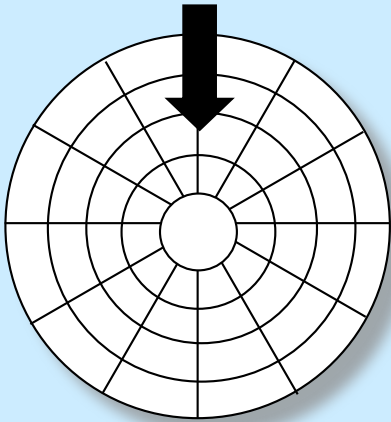
# Disk Structure: Top View of Single Platter



**Surface organized into tracks**

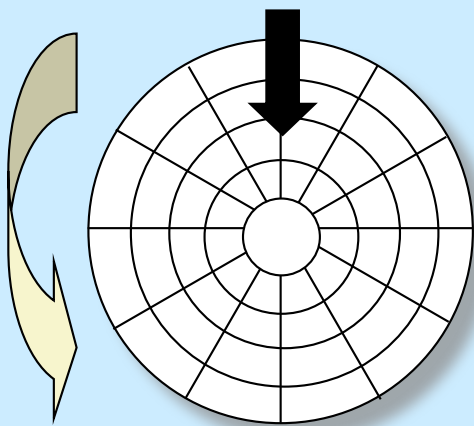
**Tracks divided into sectors**

# Disk Access



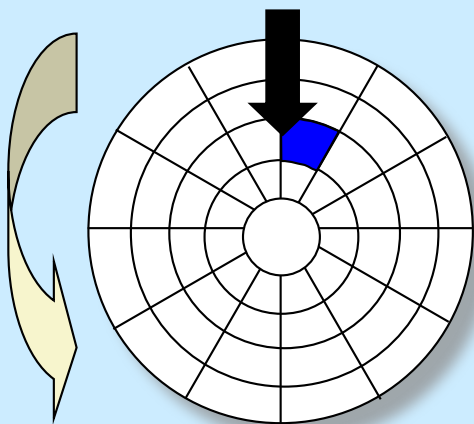
**Head in position above a track**

# Disk Access



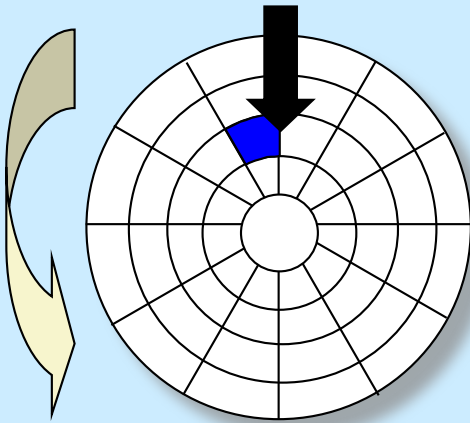
**Rotation is counter-clockwise**

# Disk Access – Read



**About to read blue sector**

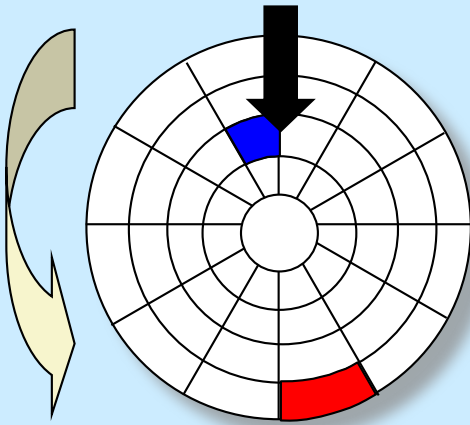
# Disk Access – Read



After **BLUE**  
read

**After reading blue sector**

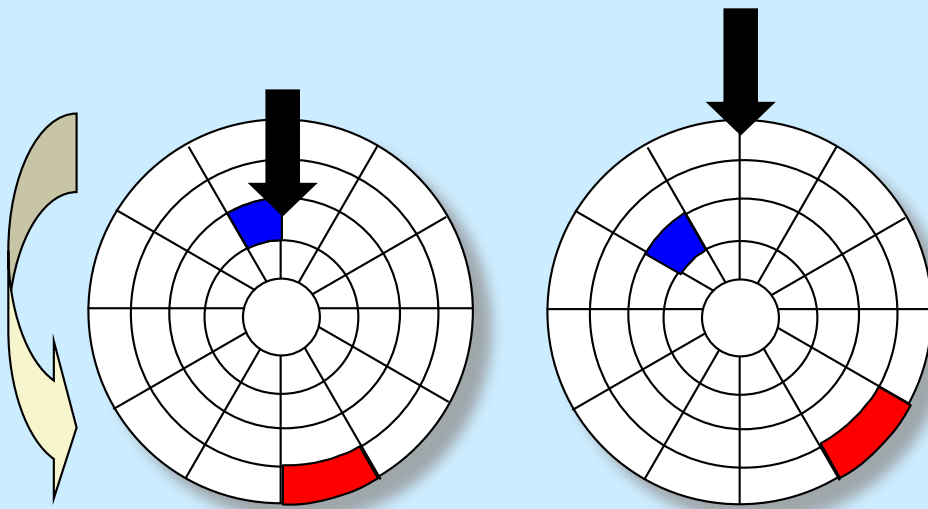
# Disk Access – Read



After **BLUE**  
read

**Red request scheduled next**

# Disk Access – Seek

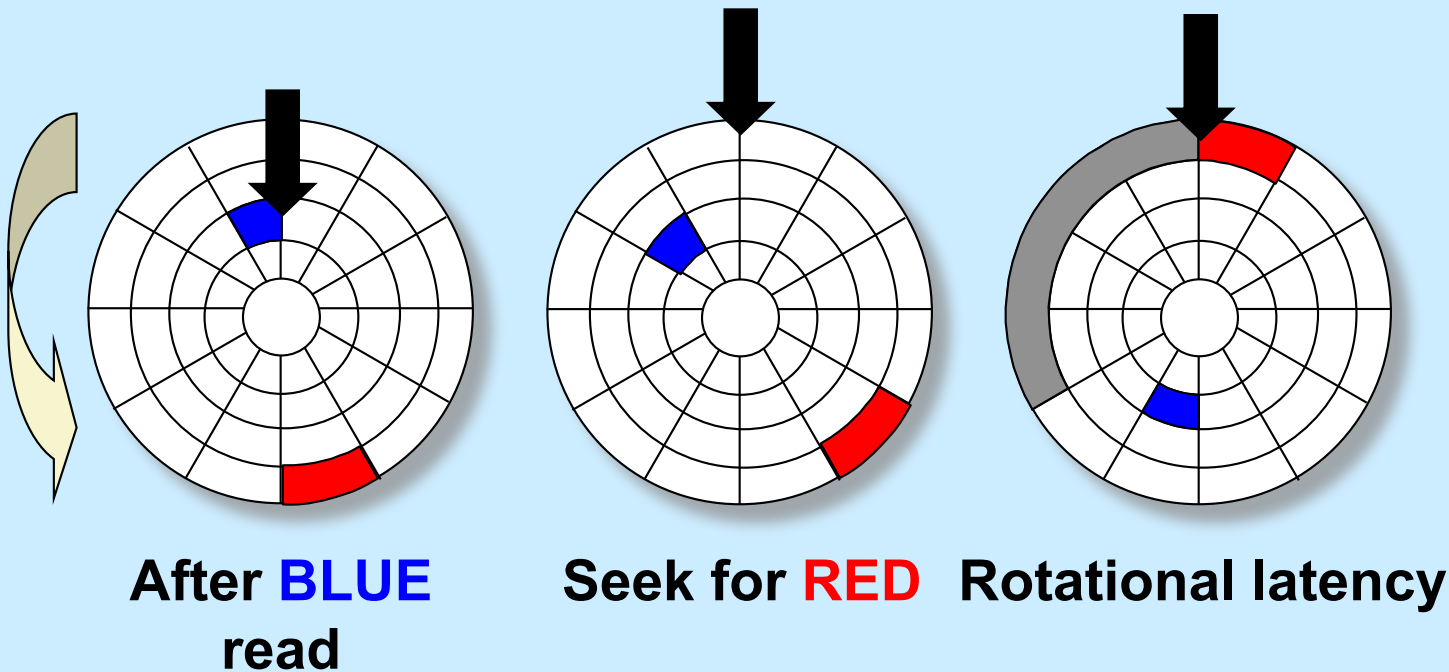


After **BLUE**  
read

Seek for **RED**

**Seek to red's track**

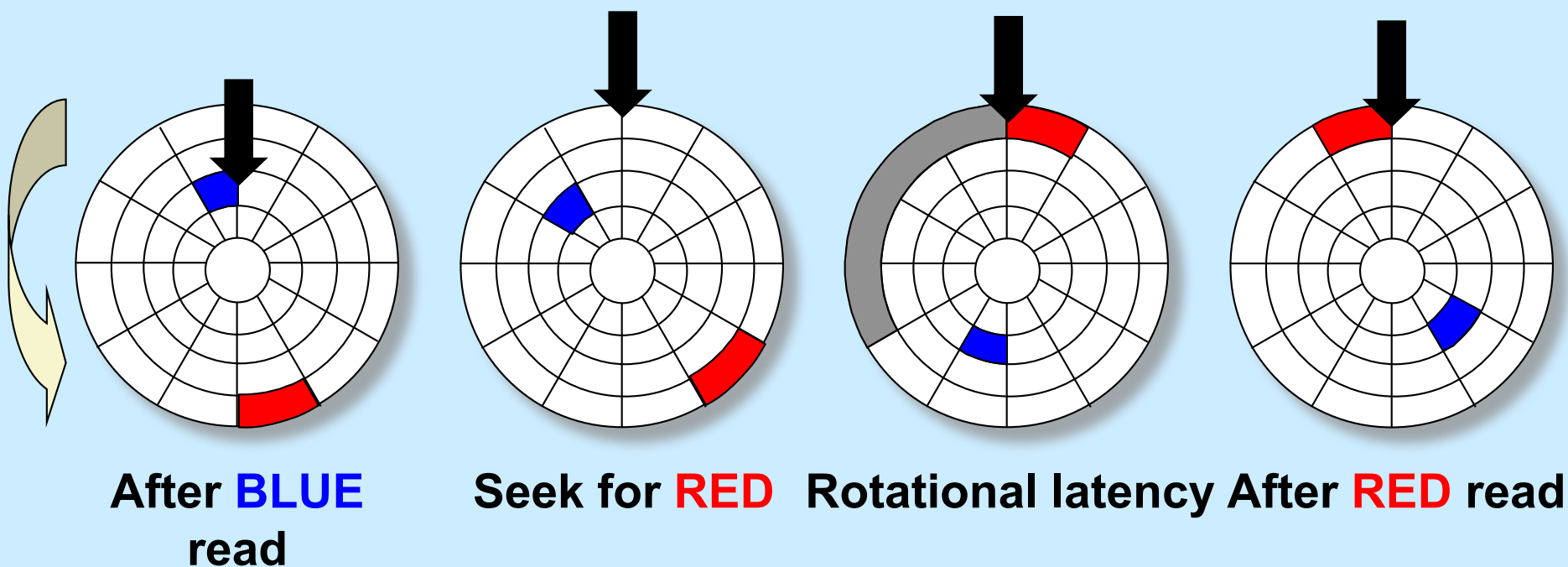
# Disk Access – Rotational Latency



**Wait for red sector to rotate around**

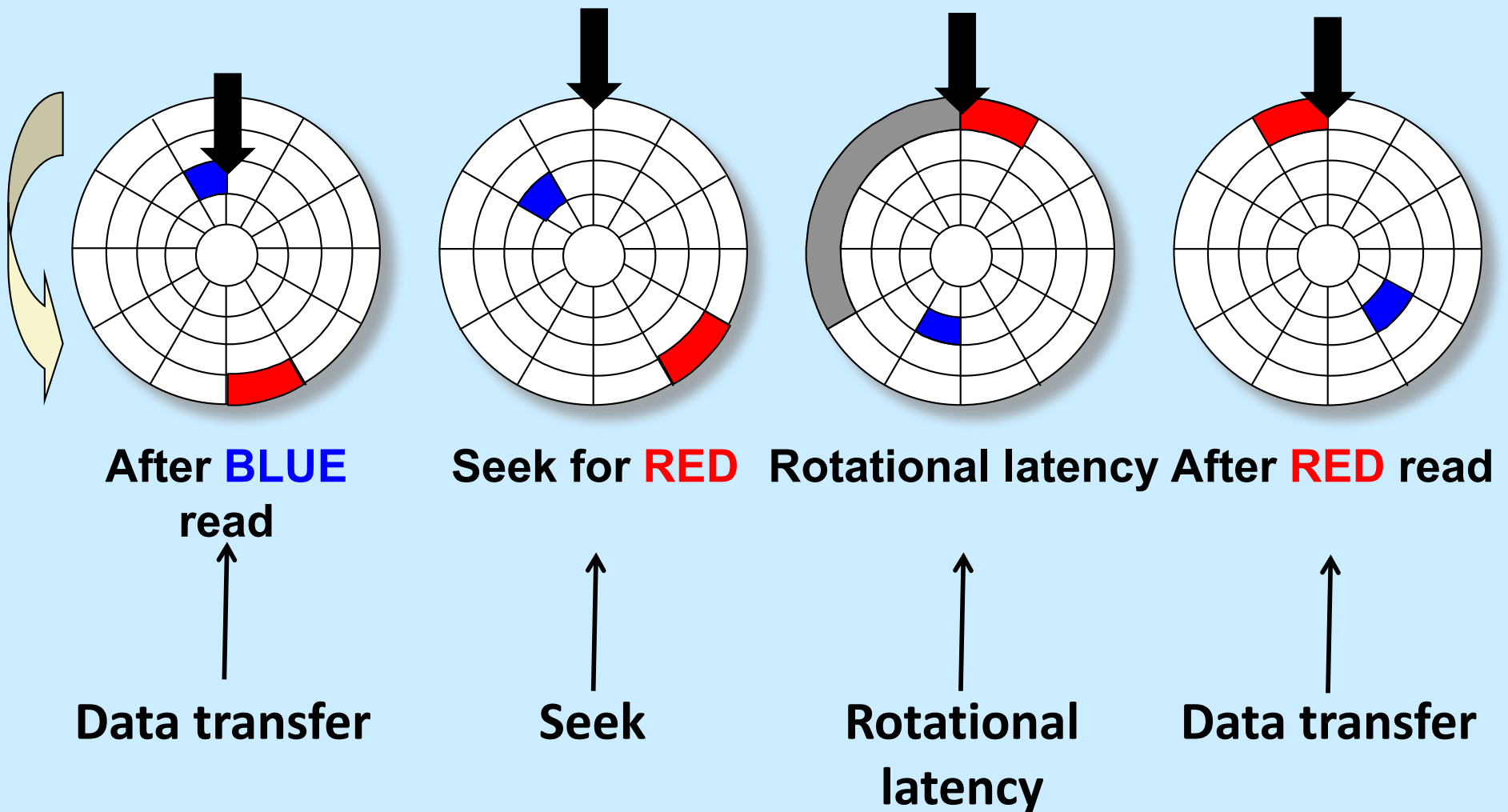


# Disk Access – Read



**Complete read of red**

# Disk Access – Service Time Components



# Disk Access Time

- Average time to access some target sector approximated 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
  - typical rotation speed  $R = 7200$  RPM
  - $T_{\text{avg rotation}} = 1/2 \times 1/R \times 60 \text{ sec/1 min}$
- **Transfer time** ( $T_{\text{avg transfer}}$ )
  - time to read the bits in the target sector
  - $T_{\text{avg transfer}} = 1/R \times 1/(\text{avg \# sectors/track}) \times 60 \text{ secs/1 min}$

# Disk Access Time Example

- **Given:**
  - rotational rate = 7,200 RPM
  - average seek time = 9 ms
  - avg # sectors/track = 600
- **Derived:**
  - $T_{\text{avg rotation}} = 1/2 \times (60 \text{ secs}/7200 \text{ RPM}) \times 1000 \text{ ms/sec} = 4 \text{ ms}$
  - $T_{\text{avg transfer}} = 60/7200 \text{ RPM} \times 1/600 \text{ sects/track} \times 1000 \text{ ms/sec} = 0.014 \text{ ms}$
  - $T_{\text{access}} = 9 \text{ ms} + 4 \text{ ms} + 0.014 \text{ ms}$
- **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

# Quiz 1

**Assuming a 5-inch diameter disk spinning at 10,000 RPM, what is the approximate speed at which the outermost track is moving?**

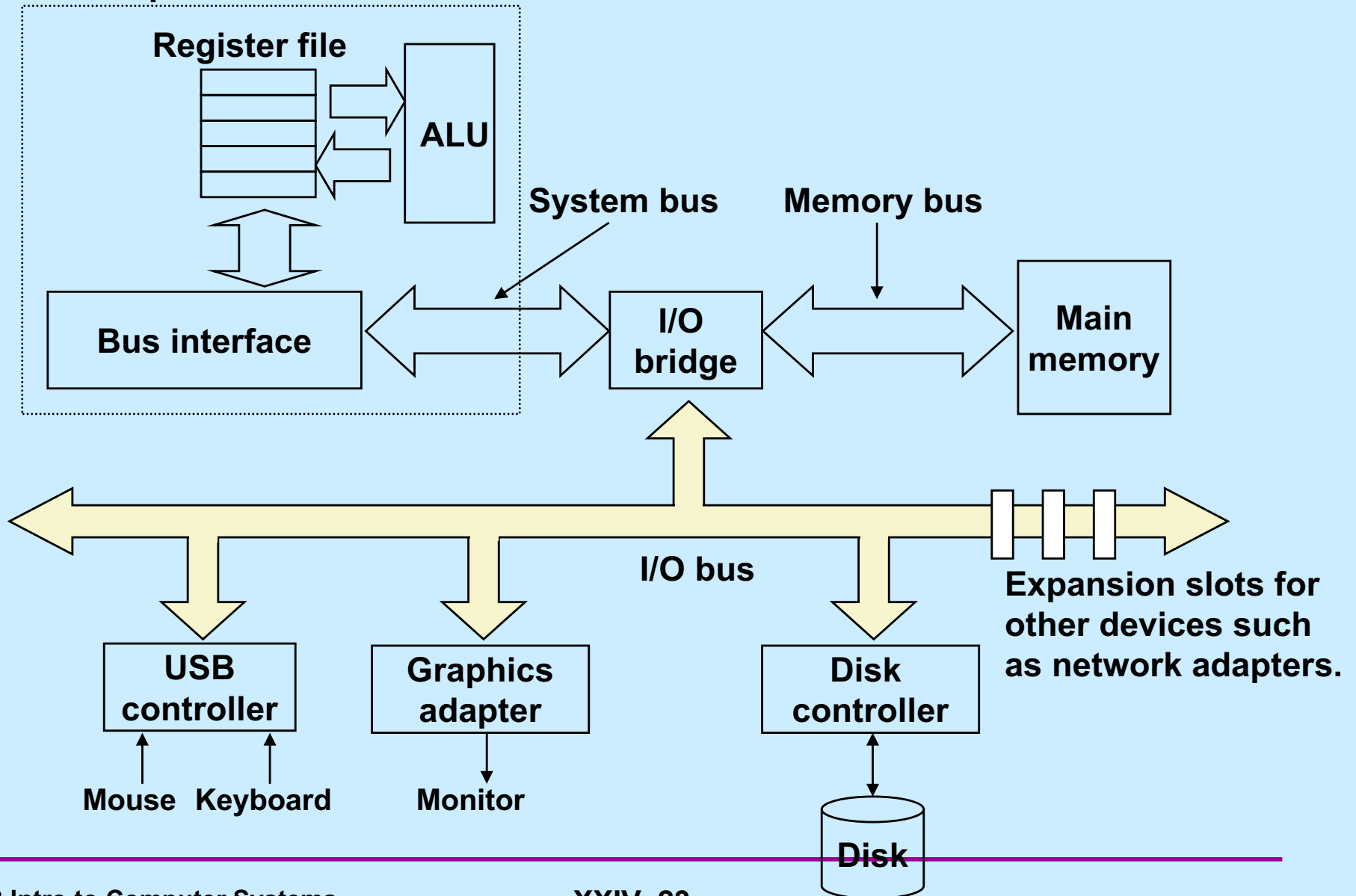
- a) faster than a speeding bullet (i.e., supersonic)**
- b) roughly the speed of a pretty fast car (250 kph/155 mph)**
- c) roughly the speed of a pretty slow car (50 mph)**
- d) roughly the speed of a world-class marathoner (13.1 mph)**

# 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”
-

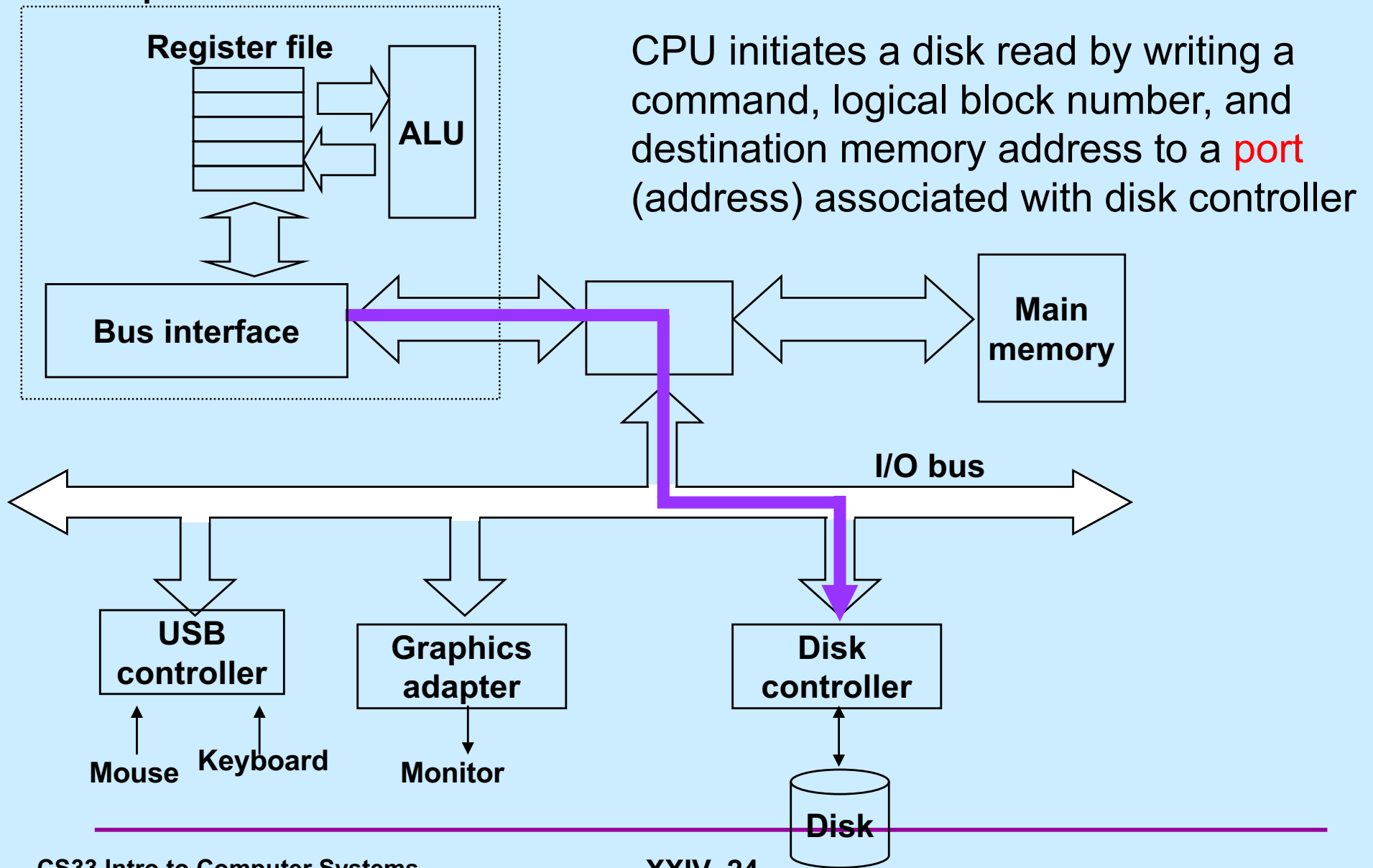
# I/O Bus

CPU chip



# Reading a Disk Sector (1)

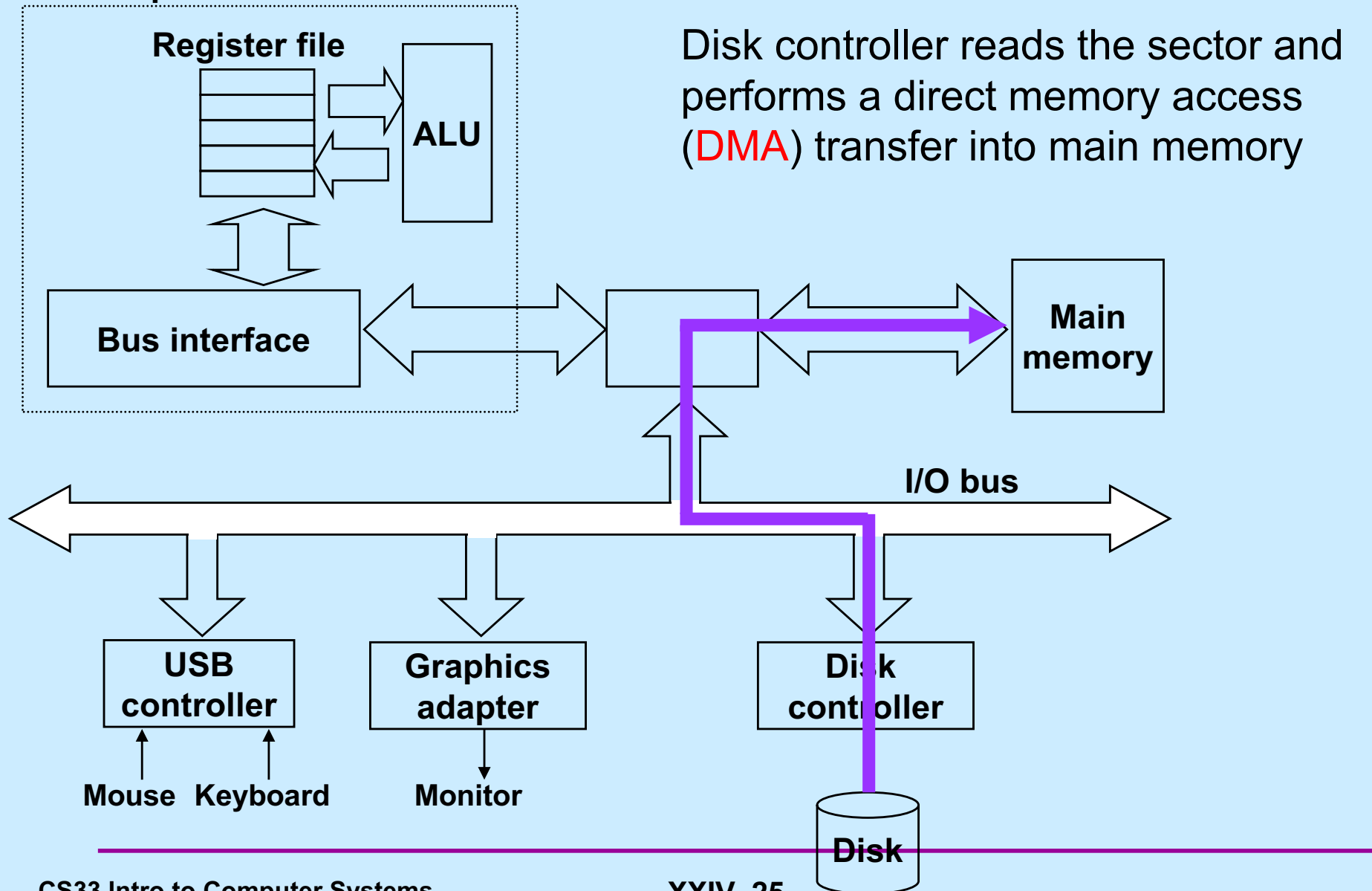
CPU chip





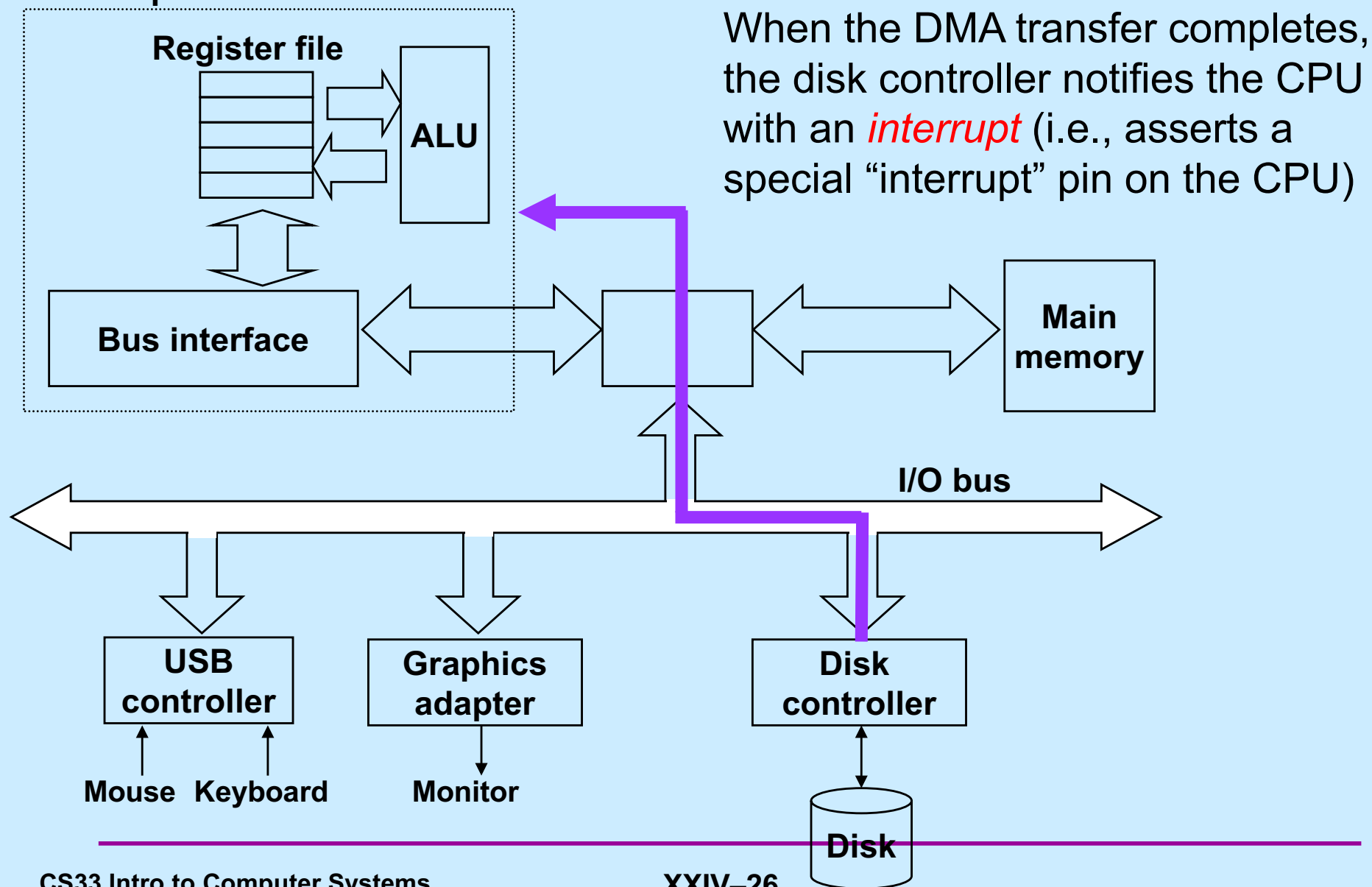
# Reading a Disk Sector (2)

CPU chip

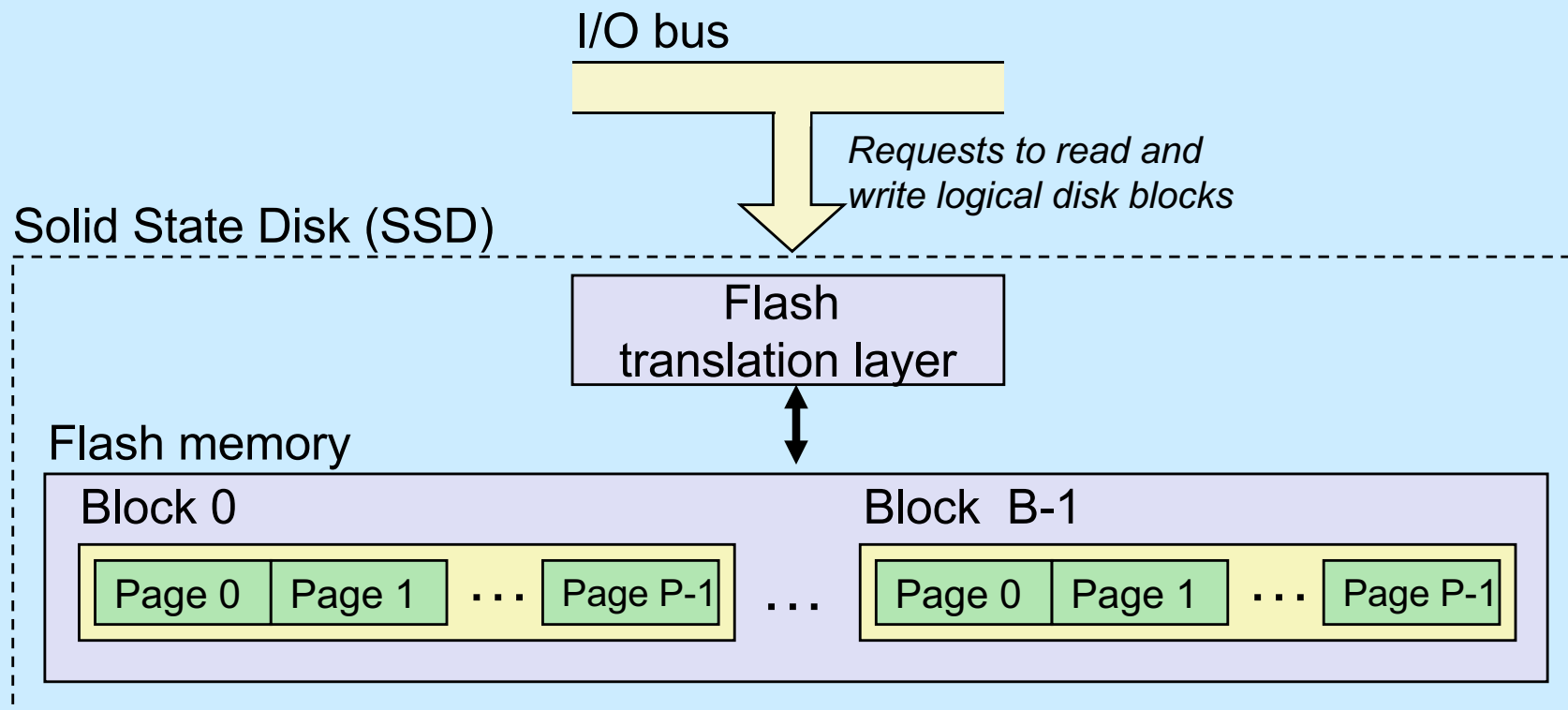


# Reading a Disk Sector (3)

CPU chip



# 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
Random read access	30 us	Random write access	300 us

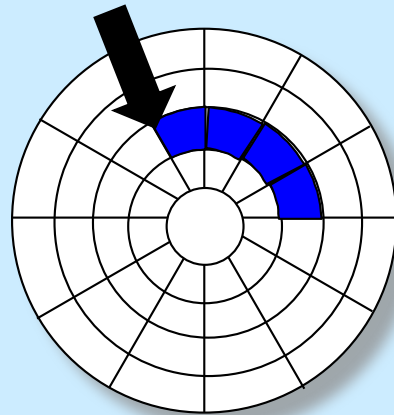
- **Why are random writes so slow?**
  - erasing a block is slow (around 1 ms)
  - modifying a page triggers a copy of all useful pages in the block
    - » find a used block (new block) and erase it
    - » write the page into the new block
    - » copy other pages from old block to the new block

# 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 ( $10^{15}$  bytes) of random writes before they wear out
  - in 2010, about 100 times more expensive per byte
  - in 2017, about 6 times more expensive per byte
- **Applications**
  - smart phones, laptops
  - Apple “Fusion” drives

# Reading a File on a Rotating Disk

- **Suppose the data of a file are stored on consecutive disk sectors on one track**
  - **this is the best possible scenario for reading data quickly**
    - » **single seek required**
    - » **single rotational delay**
    - » **all sectors read in a single scan**

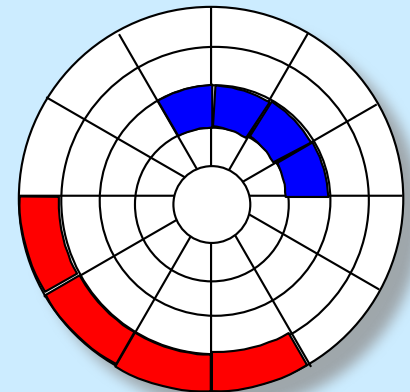


# Quiz 2

We have two files on the same (rotating) disk. The first file's data resides in consecutive sectors on one track, the second in consecutive sectors on another track. It takes a total of  $t$  seconds to read all of the first file then all of the second file.

Now suppose the files are read concurrently, perhaps a sector of the first, then a sector of the second, then the first, then the second, etc. Compared to reading them sequentially, this will take

- a) less time
- b) about the same amount of time
- c) more time



# Quiz 3

You've replaced the rotating disk on your computer with a solid-state disk. The data of the two files are again each in consecutive locations. Suppose it still takes a total of  $t$  seconds to read the first file then the second. Suppose it took  $u$  seconds to read the two files concurrently on the rotating disk. It takes  $v$  seconds to read them concurrently on the SSD.

- a)  $v < u$  (faster on the SSD)
- b)  $v \approx u$  (about the same)
- c)  $v > u$  (slower on the SSD)



# 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

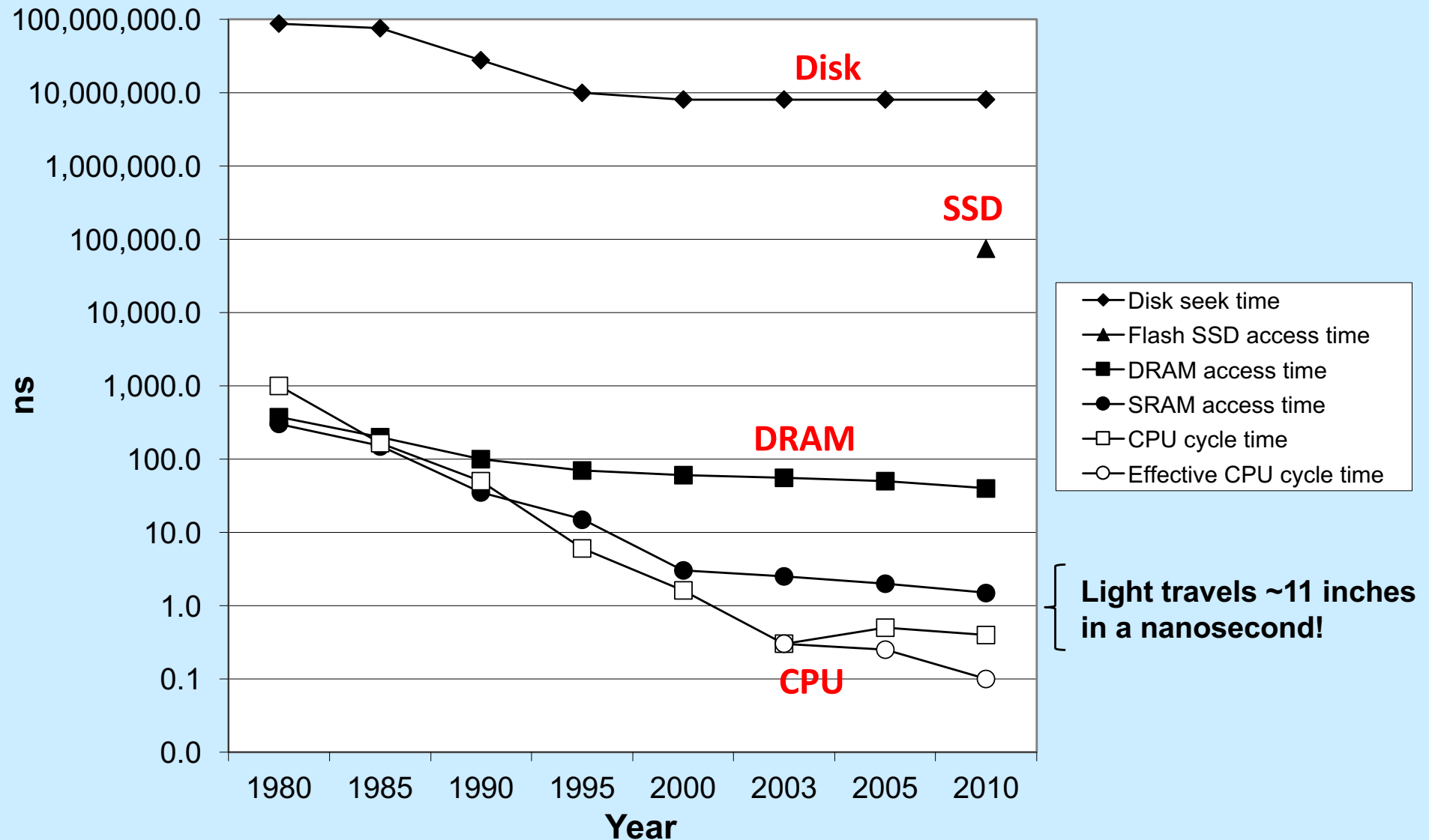
Inflection point in computer history  
when designers hit the “Power Wall”



	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

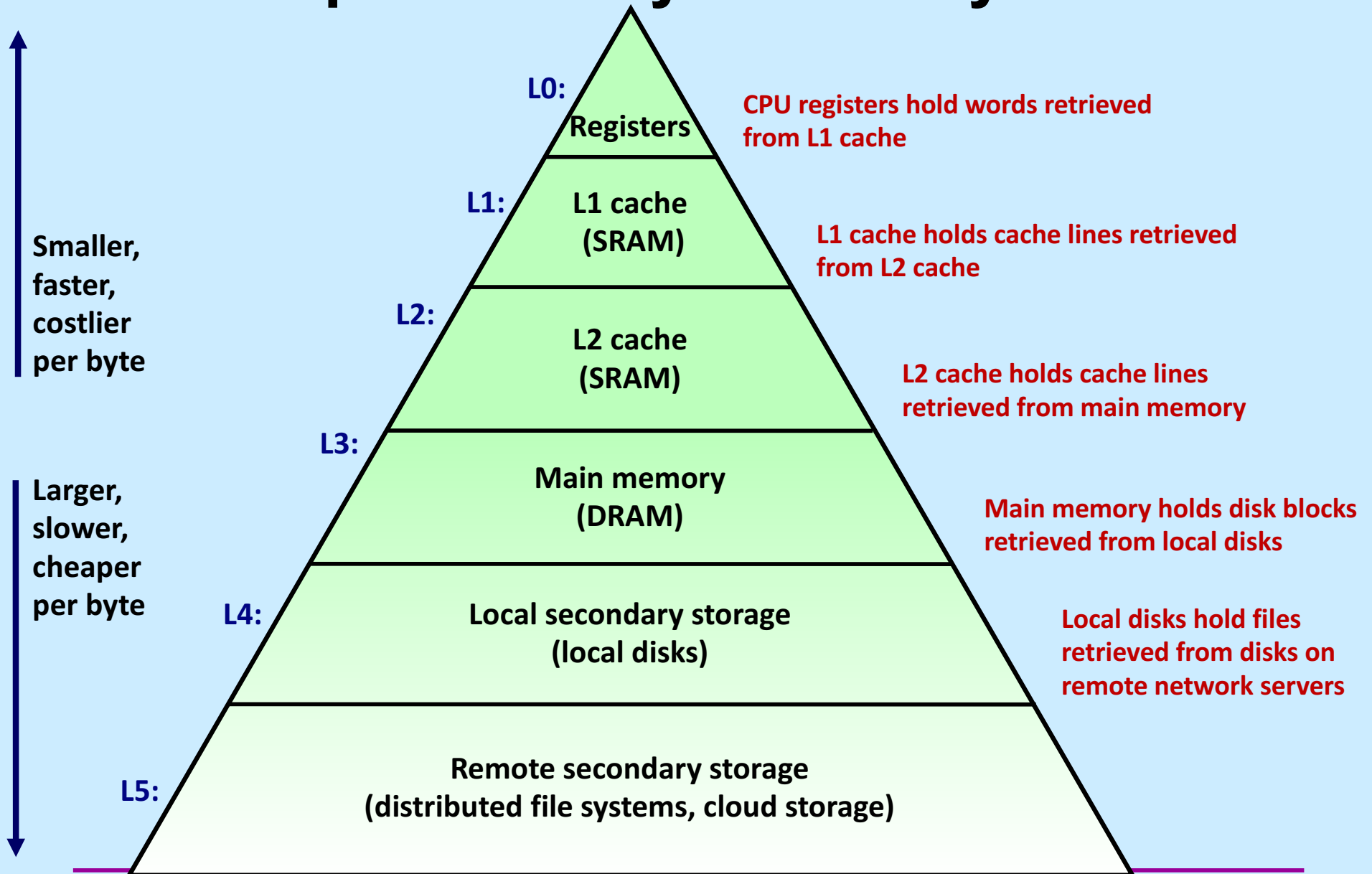
The gap widens between DRAM, disk, and CPU speeds



# 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****

# An Example Memory Hierarchy



# Putting Things Into Perspective ...

- **Reading from:**
  - ... the L1 cache is like grabbing a piece of paper from your desk (3 seconds)
  - ... the L2 cache is picking up a book from a nearby shelf (14 seconds)
  - ... main system memory is taking a 4-minute walk down the hall to talk to a friend
  - ... a hard drive is like leaving the building to roam the earth for one year and three months