

# CS 33

## Signals Part 3

# Reaping: Zombie Elimination

- **Shell must call `waitpid` on each child**
  - easy for a foreground child
  - what about background?

```
pid_t waitpid(pid_t pid, int *status, int options);
```

– ***pid* values:**

- < -1    any child process whose process group is |pid|
- 1      any child process
- 0      any child process whose process group is that of caller
- > 0    child process whose ID is equal to pid

– `wait(&status)` **is equivalent to** `waitpid(-1, &status, 0)`

# (continued)

```
pid_t waitpid(pid_t pid, int *status, int options);
```

– *options* are some combination of the following

» **WNOHANG**

- return immediately if no child has exited (returns 0)

» **WUNTRACED**

- also return if a child has been stopped (suspended)

» **WCONTINUED**

- also return if a child has been continued (resumed)

# When to Call `waitpid`

- **Shell reports status only when it is about to display its prompt**
  - thus sufficient to check on background jobs just before displaying prompt

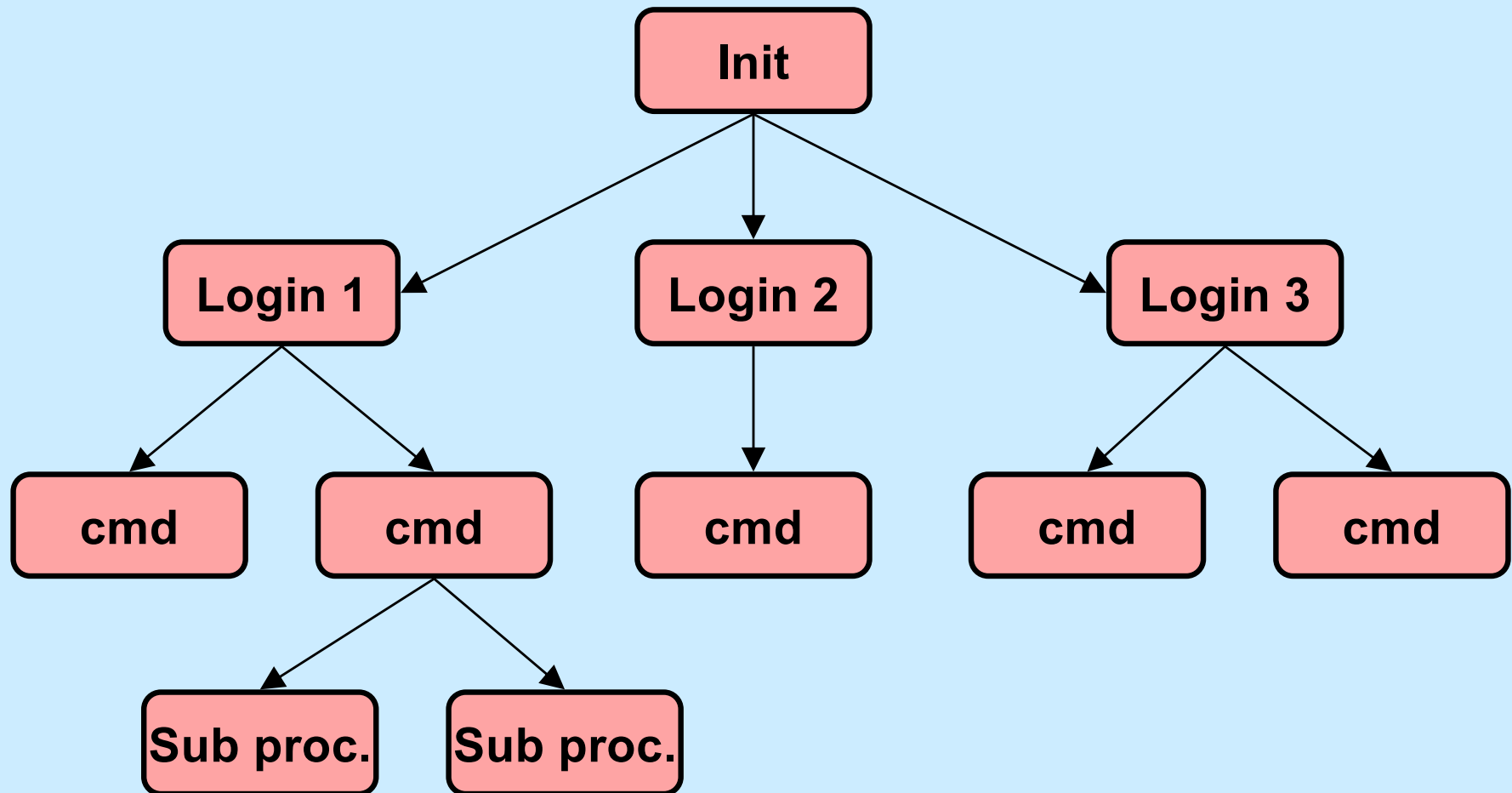
# `waitpid status`

- **WIFEXITED(\*status):** 1 if the process terminated normally and 0 otherwise
- **WEXITSTATUS(\*status):** argument to exit
- **WIFSIGNALED(\*status):** 1 if the process was terminated by a signal and 0 otherwise
- **WTERMSIG(\*status):** the signal which terminated the process if it terminated by a signal
- **WIFSTOPPED(\*status):** 1 if the process was stopped by a signal
- **WSTOPSIG(\*status):** the signal which stopped the process if it was stopped by a signal
- **WIFCONTINUED(\*status):** 1 if the process was resumed by SIGCONT and 0 otherwise

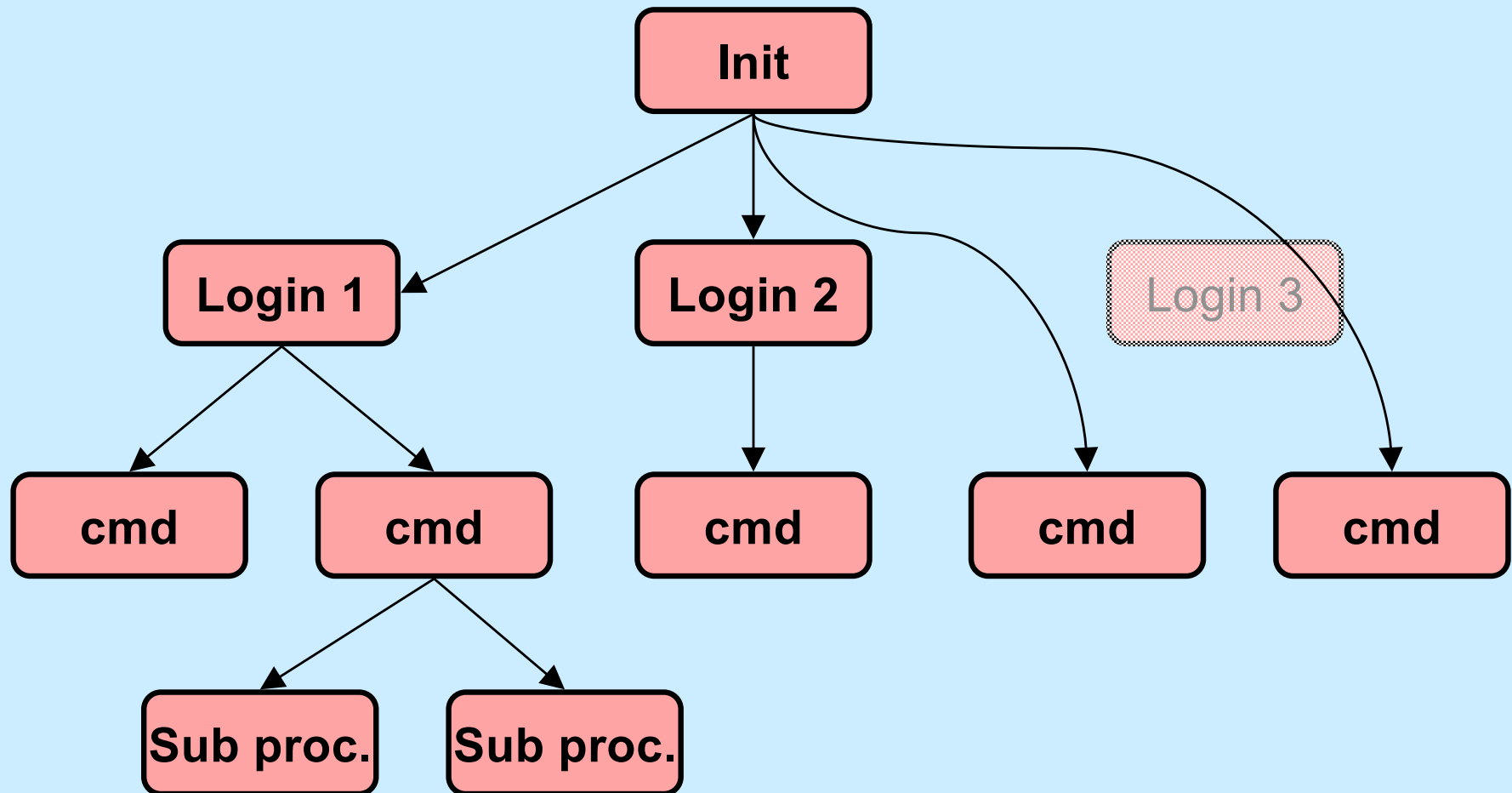
# Example (in Shell)

```
int wret, wstatus;
while ((wret = waitpid(-1, &wstatus, WNOHANG|WUNTRACED)) > 0) {
    // examine all children who've terminated or stopped
    if (WIFEXITED(wstatus)) {
        // terminated normally
        ...
    }
    if (WIFSIGNALED(wstatus)) {
        // terminated by a signal
        ...
    }
    if (WIFSTOPPED(wstatus)) {
        // stopped
        ...
    }
}
```

# Process Relationships (1)

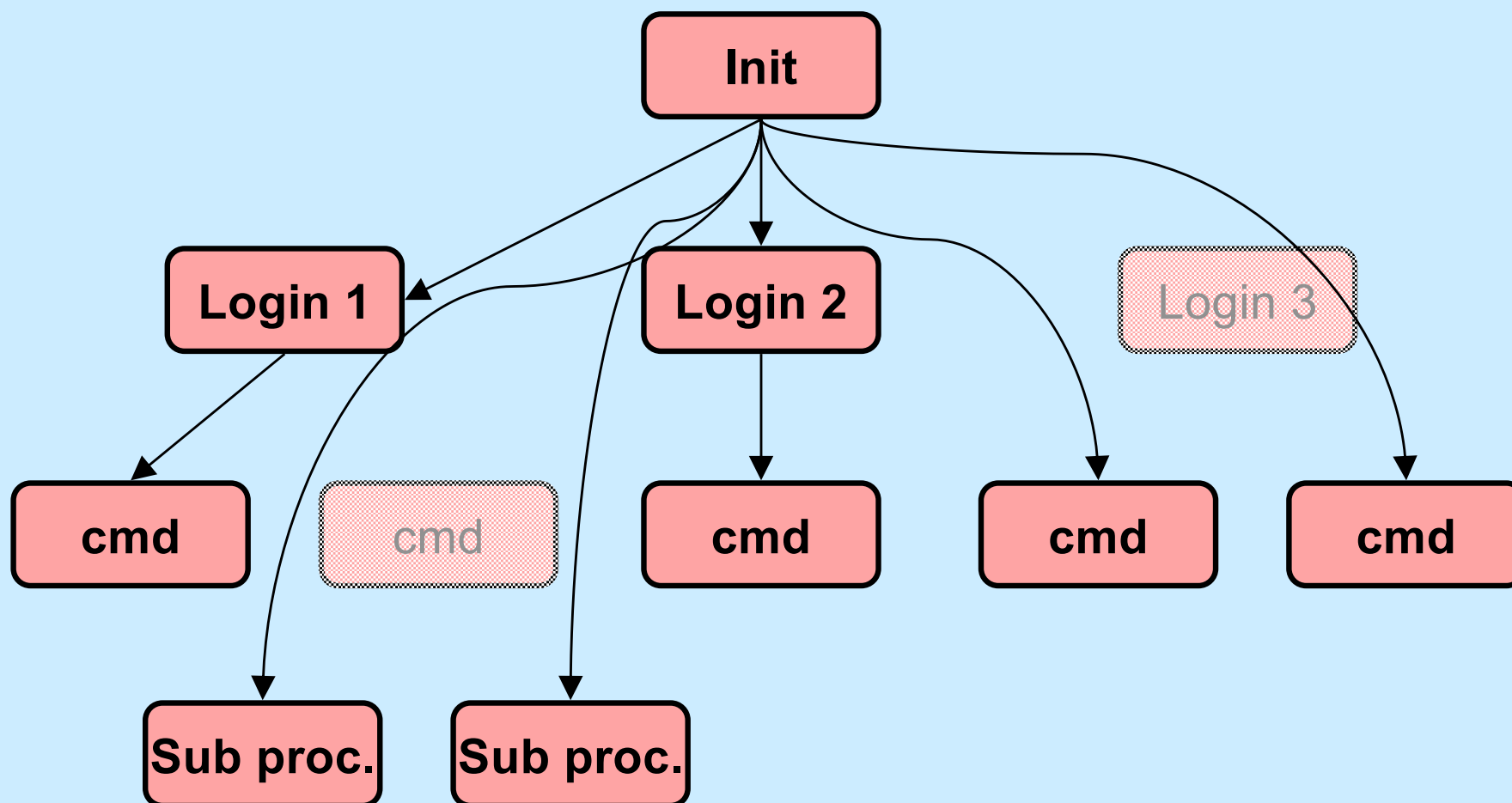


# Process Relationships (2)





# Process Relationships (3)



# Signals, Fork, and Exec

```
// set up signal handlers ...
if (fork() == 0) {
    // what happens if child gets signal?
    ...
    signal(SIGINT, SIG_IGN);
    signal(SIGFPE, handler);
    signal(SIGQUIT, SIG_DFL);
    execv("new prog", argv, NULL);
    // what happens if SIGINT, SIGFPE,
    // or SIGQUIT occur?
}
```

# Signals and System Calls

- **What happens if a signal occurs while a process is doing a system call?**
  - handler not invoked until just before system call returns to user
    - » system call might terminate early because of signal
  - system call completes
  - signal handler is invoked
  - user code resumed as if the system call has just returned

# Signals and Lengthy System Calls

- **Some system calls take a long time**
  - large I/O transfer
    - » multi-gigabyte read or write request probably done as a sequence of smaller pieces
  - a long wait is required
    - » a read from the keyboard requires waiting for someone to type something
- **If signal arrives in the midst of lengthy system call, handler invoked:**
  - after current piece is completed
  - after cancelling wait

# Interrupted System Calls

- **What if a signal is handled before the system call completes?**
  - **invoke handler, then return from system call prematurely**
    - **if one or more pieces were completed, return total number of bytes transferred**
    - **otherwise return “interrupted” error**

# Summary: Signals Occurring During System Calls

- **Either**
  - wait for system call to finish, then invoke handler
  - or
  - stop system call early, then invoke handler
    - » EINTR error if nothing had been done yet
    - » return partial results if it was underway

# Interrupted System Calls: Lengthy Case

```
char buf[BSIZE];
fillbuf(buf);
long remaining = BSIZE;
char *bptr = buf;
while (1) {
    long num_xfrd = write(fd,
        bptr, remaining);
    if (num_xfrd == -1) {
        if (errno == EINTR) {
            // interrupted early
            continue;
        }
        perror("big trouble");
        exit(1);
    }
    if (num_xfrd < remaining) {
        /* interrupted after the
           first step */
        remaining -= num_xfrd;
        bptr += num_xfrd;
        continue;
    }
    // success!
    break;
}
```

# Asynchronous Signals (1)

```
main( ) {  
    void handler(int);  
    signal(SIGINT, handler);  
  
    ... /* long-running buggy code */  
  
}  
  
void handler(int sig) {  
    ... /* clean up */  
    exit(1);  
}
```



# Asynchronous Signals (2)

```
computation_state_t state;    long_running_procedure( ) {
                                while (a_long_time) {
main( ) {                       update_state(&state);
    void handler(int);         compute_more( );
                                }
                                }
    signal(SIGINT, handler);
                                }
    long_running_procedure( );  void handler(int sig) {
}                                display(&state);
                                }

```

# Asynchronous Signals (3)

```
main( ) {  
    void handler(int);  
  
    signal(SIGINT, handler);  
  
    ... /* complicated program */  
  
    myputs("important message\n");  
  
    ... /* more program */  
  
}  
  
void handler(int sig) {  
    ... /* deal with signal */  
  
    myputs("equally important "  
           "message\n");  
}
```

# Asynchronous Signals (4)

```
char buf[BSIZE];
int pos;

void myputs(char *str) {
    int len = strlen(str);
    for (int i=0; i<len; i++, pos++) {
        buf[pos] = str[i];
        if ((buf[pos] == '\n') || (pos == BSIZE-1)) {
            write(1, buf, pos+1);
            pos = -1;
        }
    }
}
```

# Async-Signal Safety

- Which library functions are safe to use within signal handlers?

- abort	- dup2	- getppid	- readlink	- sigemptyset	- tcgetpgrp
- accept	- execl	- getsockname	- recv	- sigfillset	- tcsendbreak
- access	- execve	- getsockopt	- recvfrom	- sigismember	- tcsetattr
- aio_error	- _exit	- getuid	- recvmsg	- signal	- tcsetpgrp
- aio_return	- fchmod	- kill	- rename	- sigpause	- time
- aio_suspend	- fchown	- link	- rmdir	- sigpending	- timer_getoverrun
- alarm	- fcntl	- listen	- select	- sigprocmask	- timer_gettime
- bind	- fdatsync	- lseek	- sem_post	- sigqueue	- timer_settime
- cfgetispeed	- fork	- lstat	- send	- sigsuspend	- times
- cfgetospeed	- fpathconf	- mkdir	- sendmsg	- sleep	- umask
- cfsetispeed	- fstat	- mkfifo	- sendto	- sockatmark	- uname
- cfsetospeed	- fsync	- open	- setgid	- socket	- unlink
- chdir	- ftruncate	- pathconf	- setpgid	- socketpair	- utime
- chmod	- getegid	- pause	- setsid	- stat	- wait
- chown	- geteuid	- pipe	- setsockopt	- symlink	- waitpid
- clock_gettime	- getgid	- poll	- setuid	- sysconf	- write
- close	- getgroups	- posix_trace_event	- shutdown	- tcdrain	
- connect	- getpeername	- pselect	- sigaction	- tcflow	
- creat	- getpgrp	- raise	- sigaddset	- tcflush	
- dup	- getpid	- read	- sigdelset	- tcgetattr	

# Quiz 1

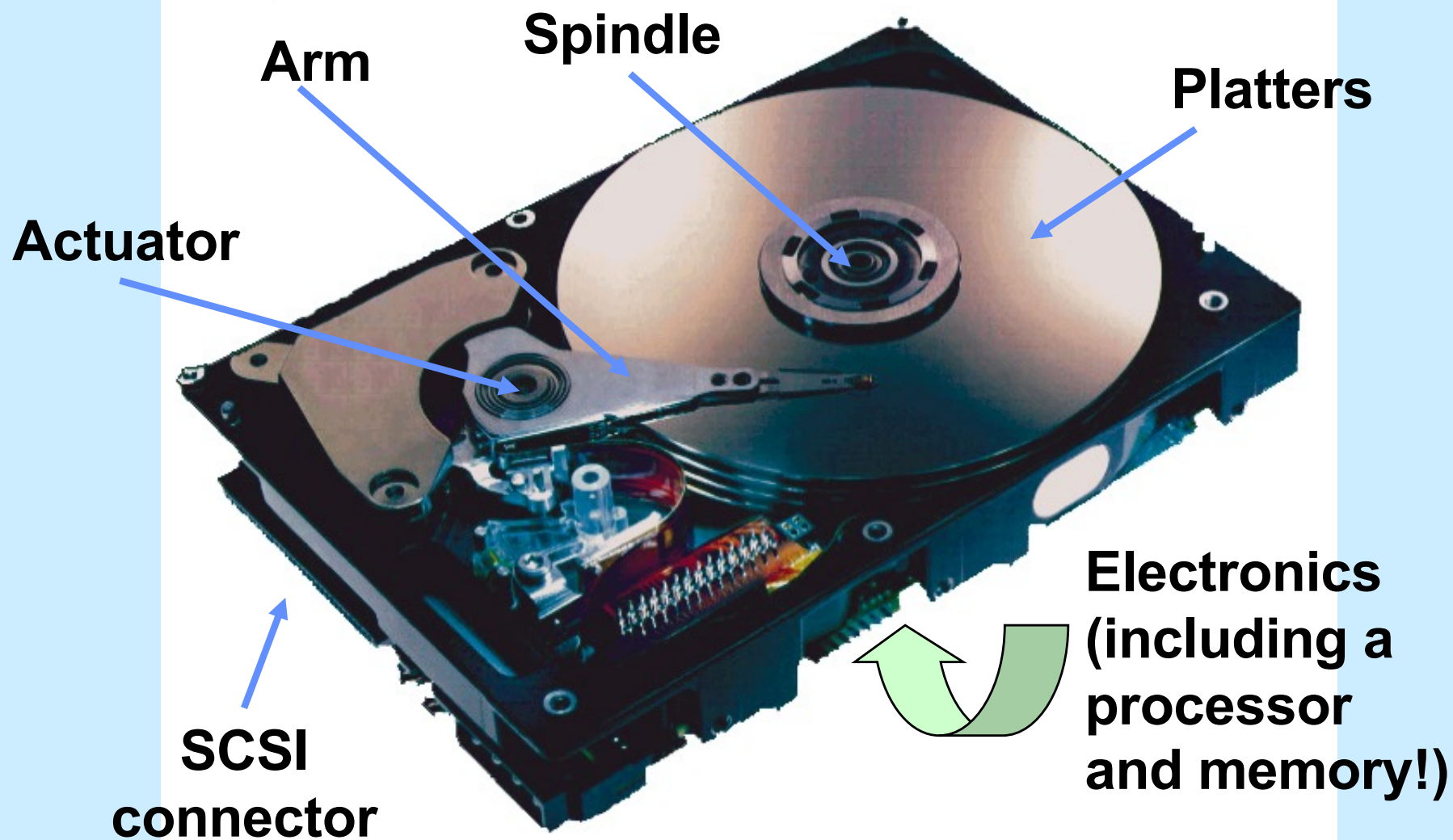
**Printf is not listed as being async-signal safe.  
Can it be implemented so that it is?**

- a) yes, but it would be so complicated, it's not done**
- b) yes, it can be easily made async-signal safe**
- c) no, it's inherently not async-signal safe**

# CS 33

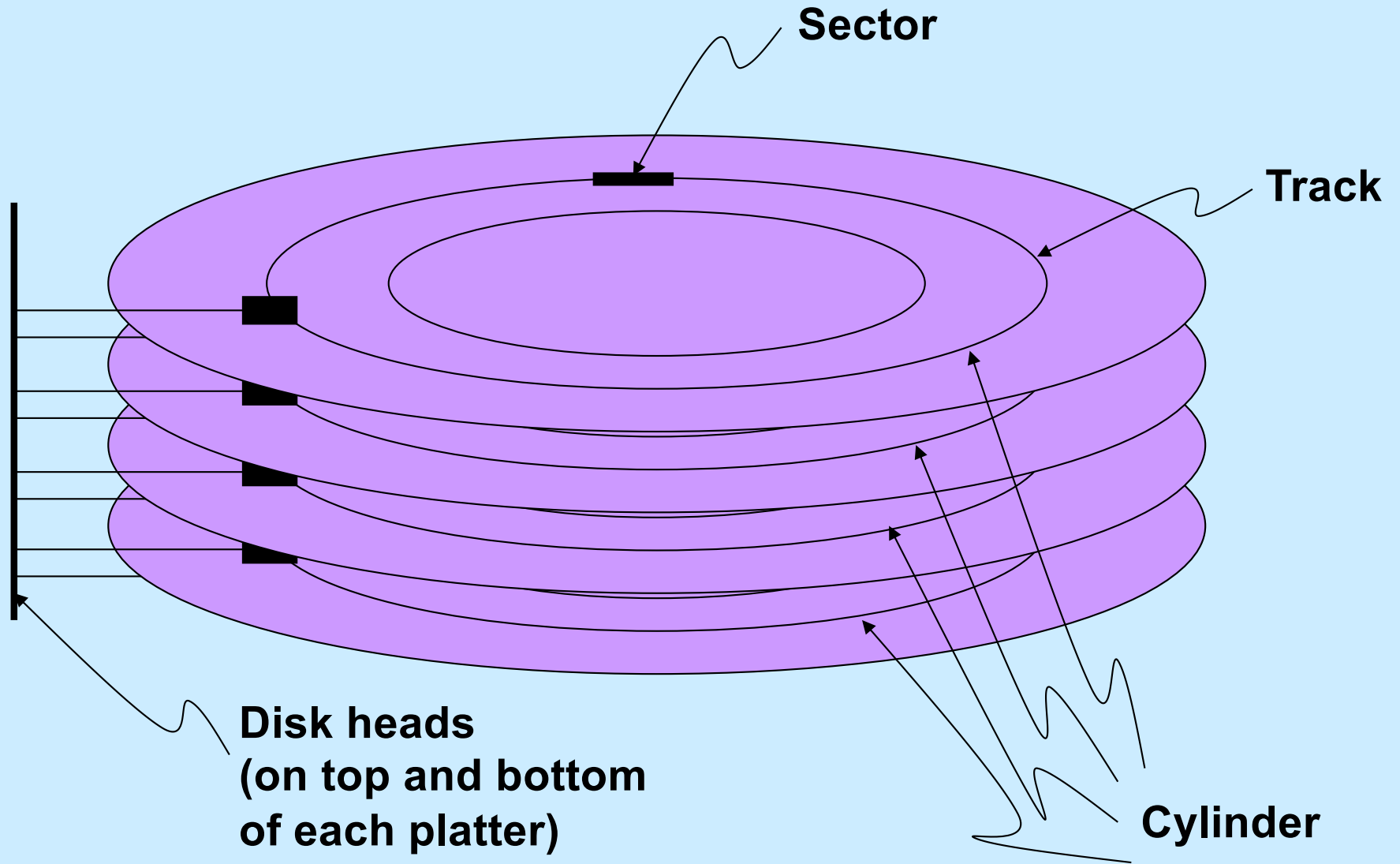
## Memory Hierarchy II

# What's Inside A Disk Drive?



*Image courtesy of Seagate Technology*

# Disk Architecture

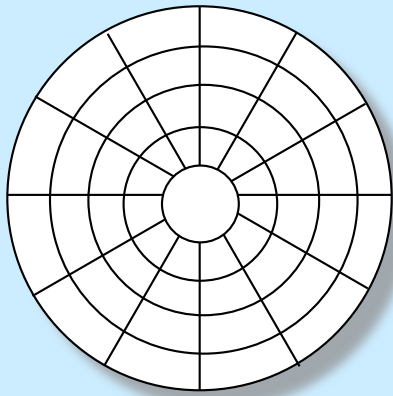




# Example Disk Drive

<b>Rotation speed</b>	<b>10,000 RPM</b>
<b>Number of surfaces</b>	<b>8</b>
<b>Sector size</b>	<b>512 bytes</b>
<b>Sectors/track</b>	<b>500-1000; 750 average</b>
<b>Tracks/surface</b>	<b>100,000</b>
<b>Storage capacity</b>	<b>307.2 billion bytes</b>
<b>Average seek time</b>	<b>4 milliseconds</b>
<b>One-track seek time</b>	<b>.2 milliseconds</b>
<b>Maximum seek time</b>	<b>10 milliseconds</b>

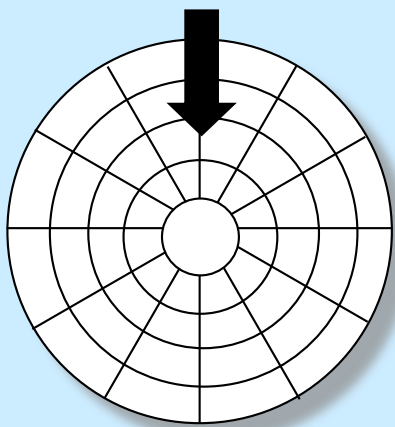
# Disk Structure: Top View of Single Platter



**Surface organized into tracks**

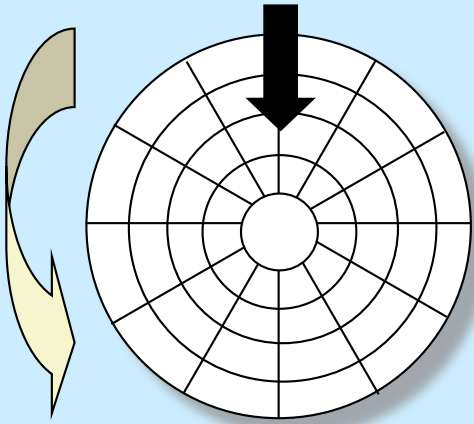
**Tracks divided into sectors**

# Disk Access



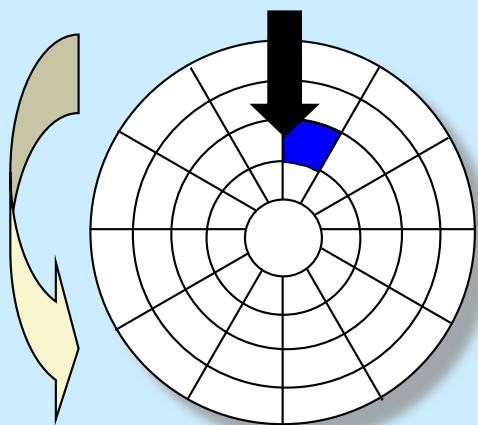
**Head in position above a track**

# Disk Access



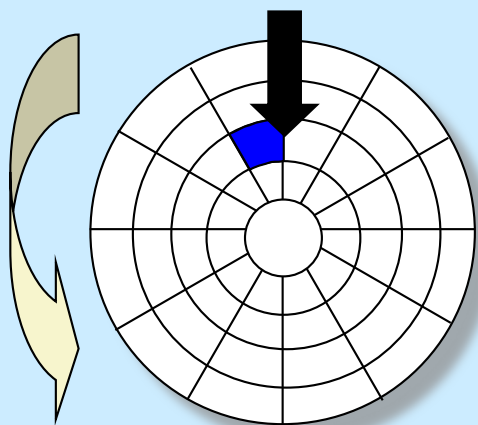
**Rotation is counter-clockwise**

# Disk Access – Read



**About to read blue sector**

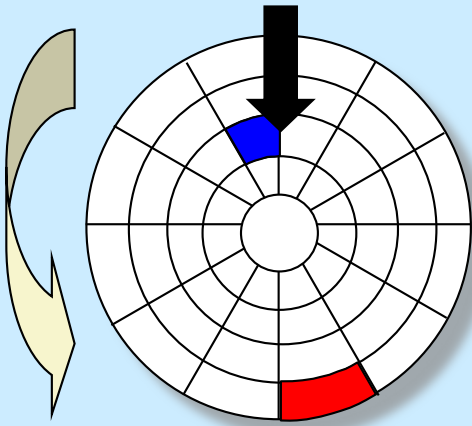
# Disk Access – Read



After **BLUE**  
read

After reading blue sector

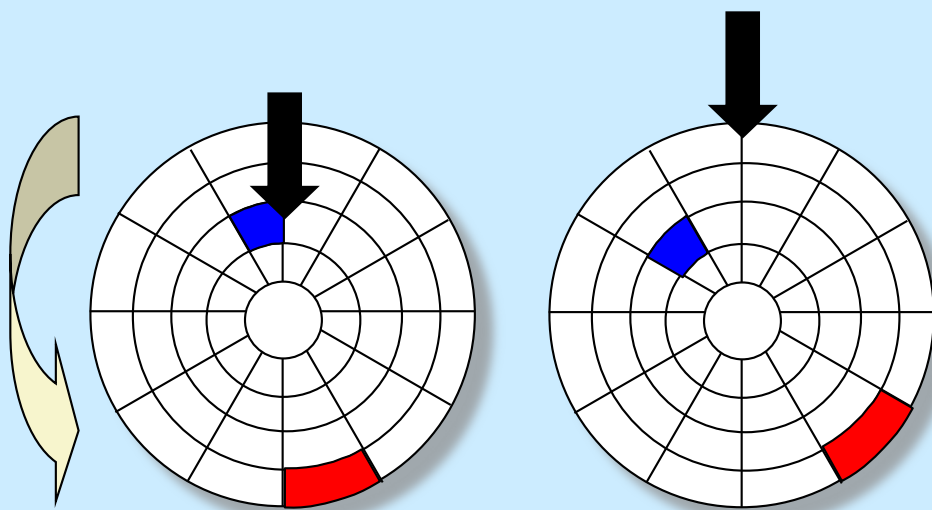
# Disk Access – Read



After **BLUE**  
read

**Red request scheduled next**

# Disk Access – Seek



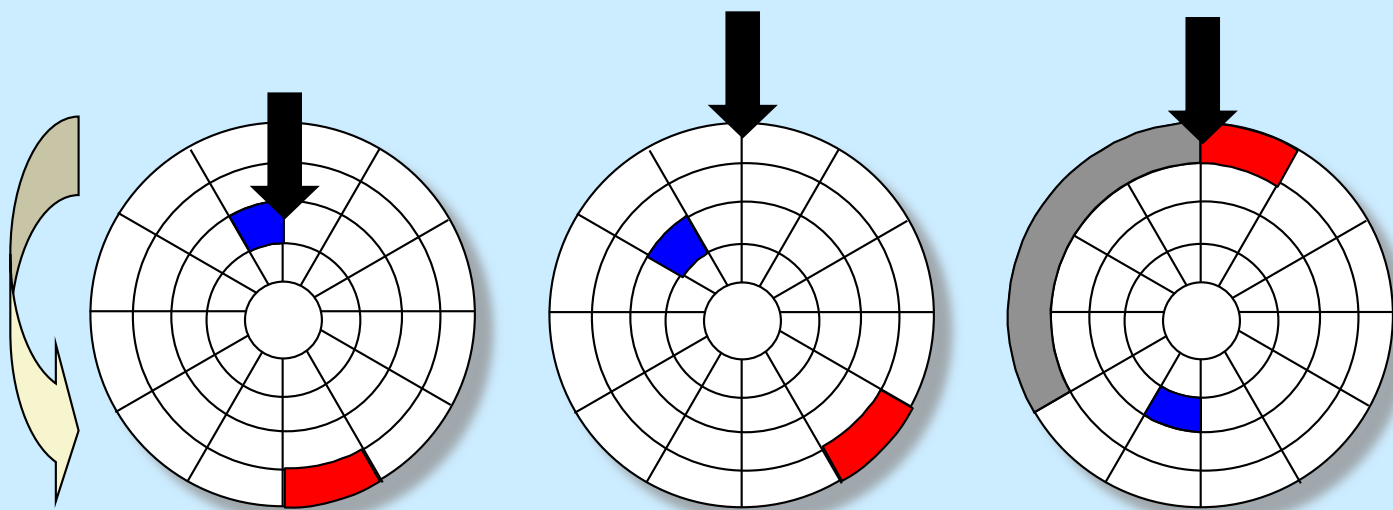
After **BLUE**  
read

Seek for **RED**

## Seek to red's track



# Disk Access – Rotational Latency

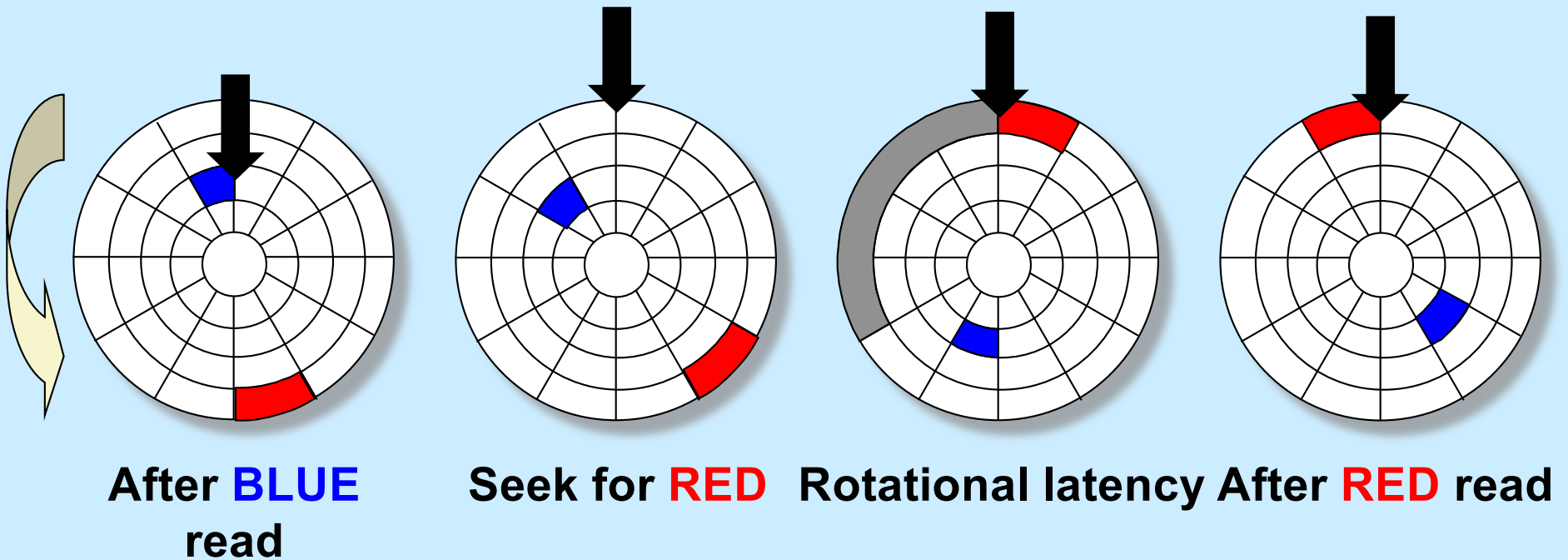


After **BLUE**  
read

Seek for **RED** Rotational latency

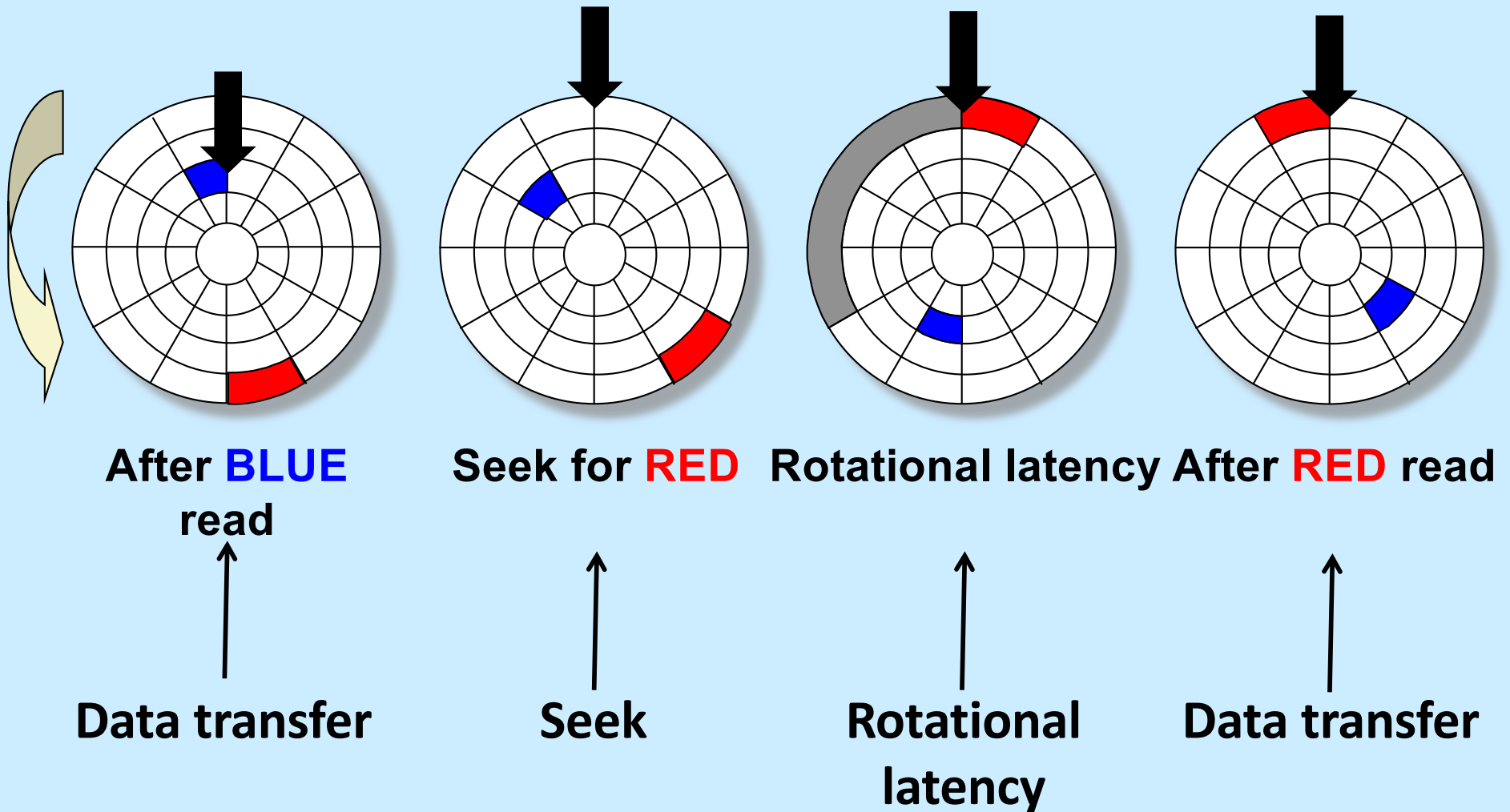
**Wait for red sector to rotate around**

# Disk Access – Read



**Complete read of red**

# Disk Access – Service Time Components



# Disk Access Time

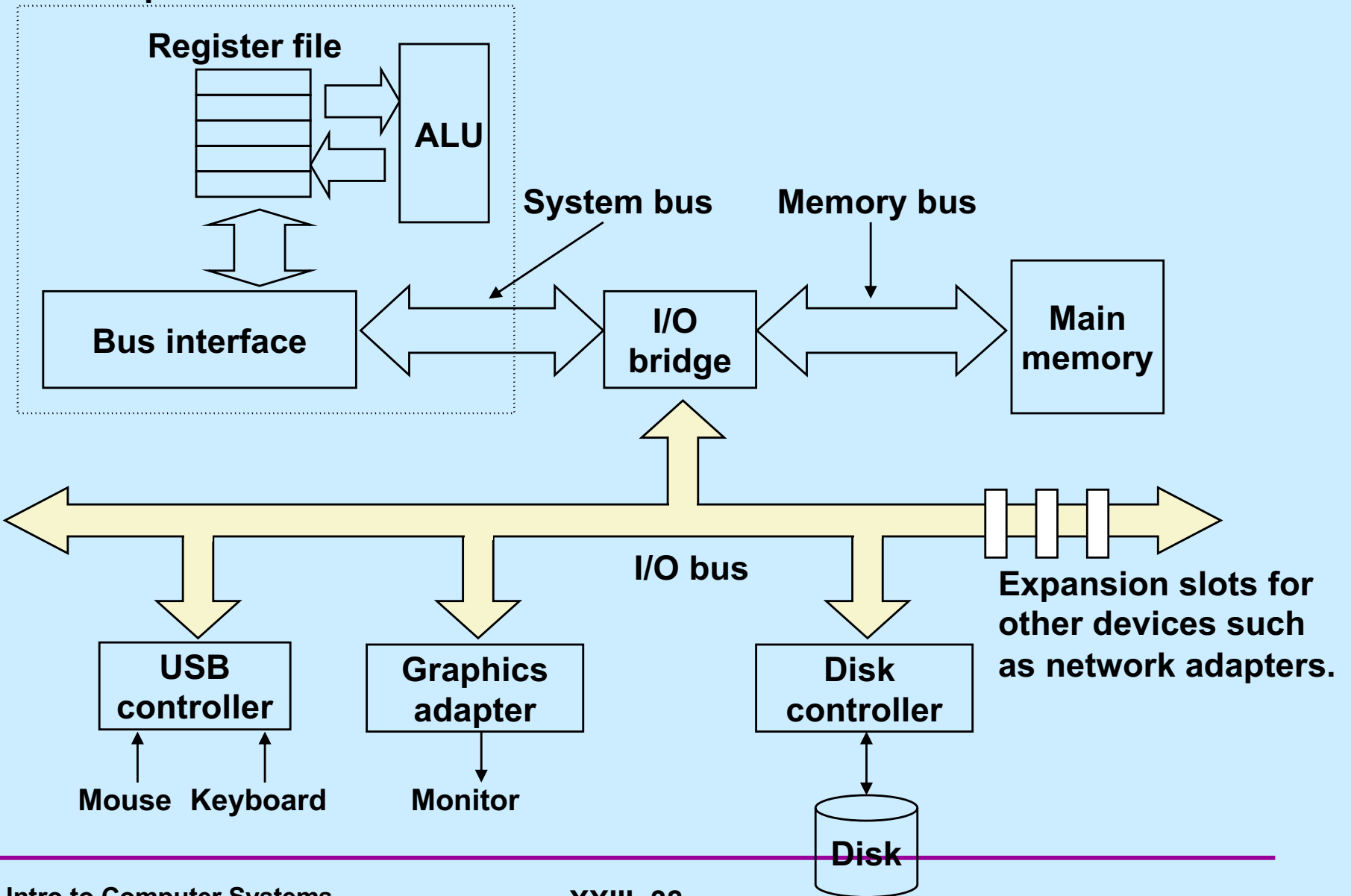
- Average time to access some target sector approximated by :
  - $T_{\text{access}} = T_{\text{avg seek}} + T_{\text{avg rotation}} + T_{\text{avg transfer}}$
- **Seek time** ( $T_{\text{avg seek}}$ )
  - time to position heads over cylinder containing target sector
  - typical  $T_{\text{avg seek}}$  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
  - $T_{\text{avg rotation}} = 1/2 \times 1/R \times 60 \text{ sec}/1 \text{ min}$
- **Transfer time** ( $T_{\text{avg transfer}}$ )
  - time to read the bits in the target sector
  - $T_{\text{avg transfer}} = 1/R \times 1/(\text{avg \# sectors/track}) \times 60 \text{ secs}/1 \text{ min}$

# 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

# I/O Bus

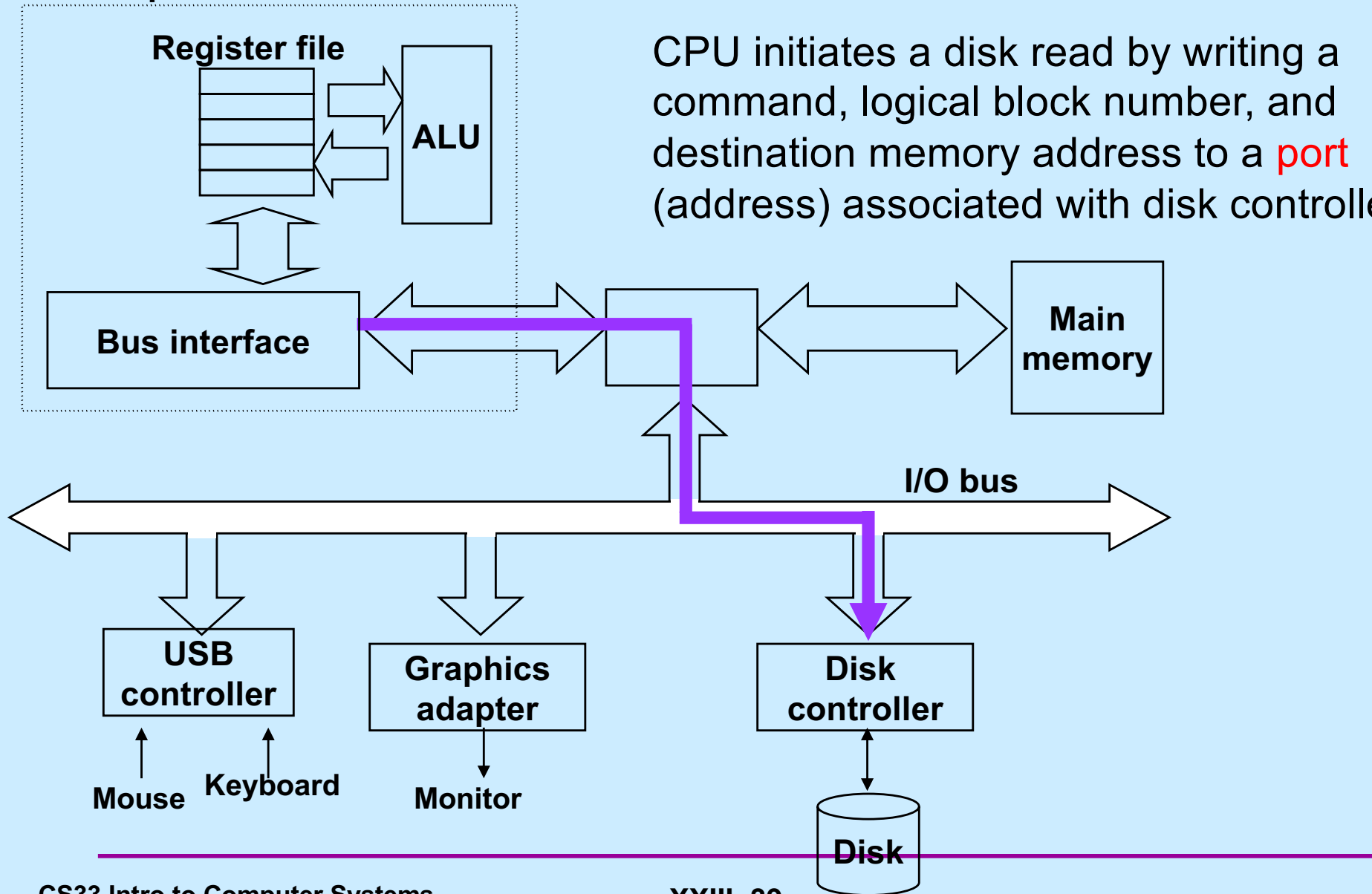
CPU chip



# Reading a Disk Sector (1)

CPU chip

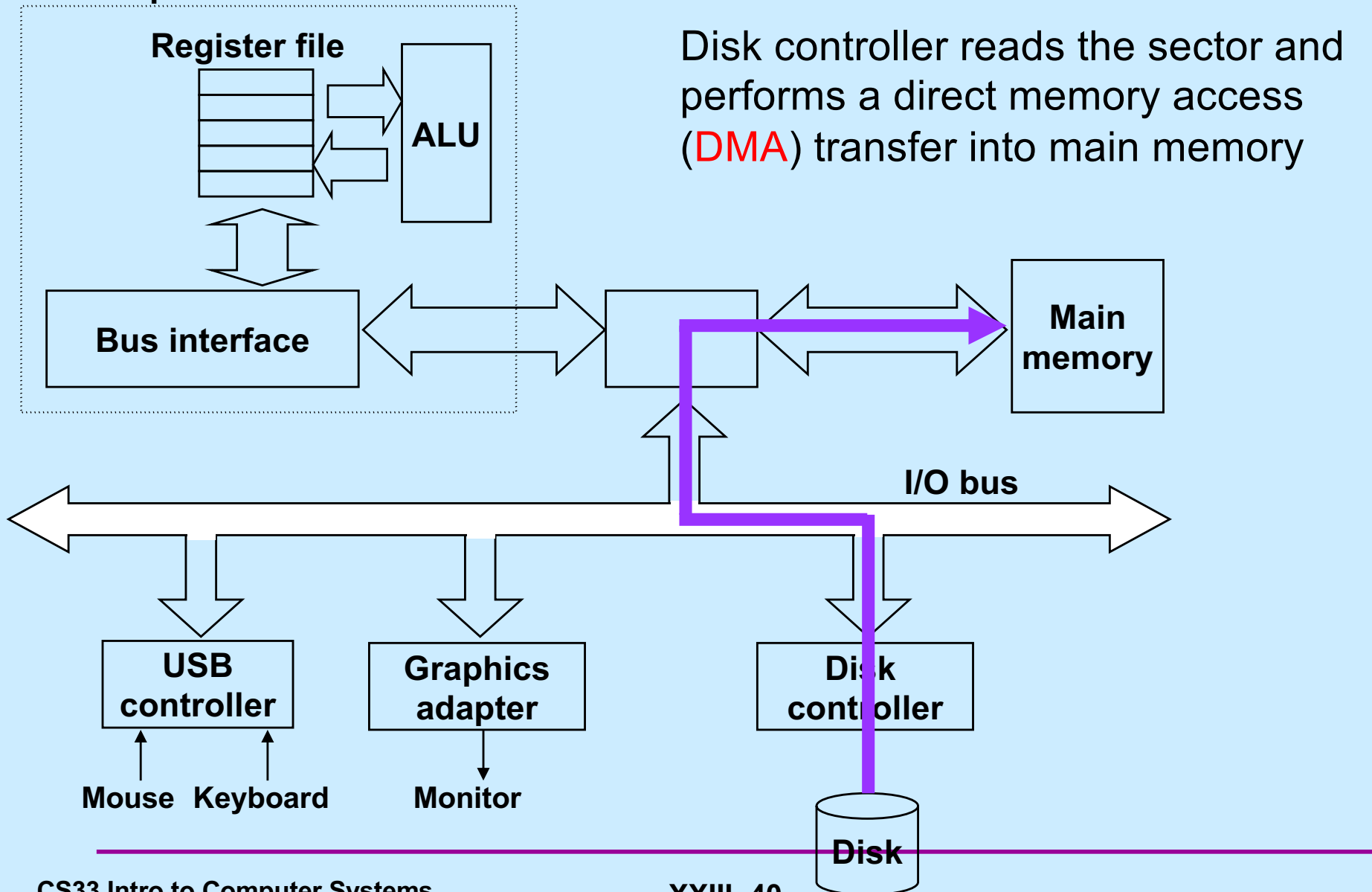
CPU initiates a disk read by writing a command, logical block number, and destination memory address to a **port** (address) associated with disk controller



# Reading a Disk Sector (2)

CPU chip

Disk controller reads the sector and performs a direct memory access (DMA) transfer into main memory

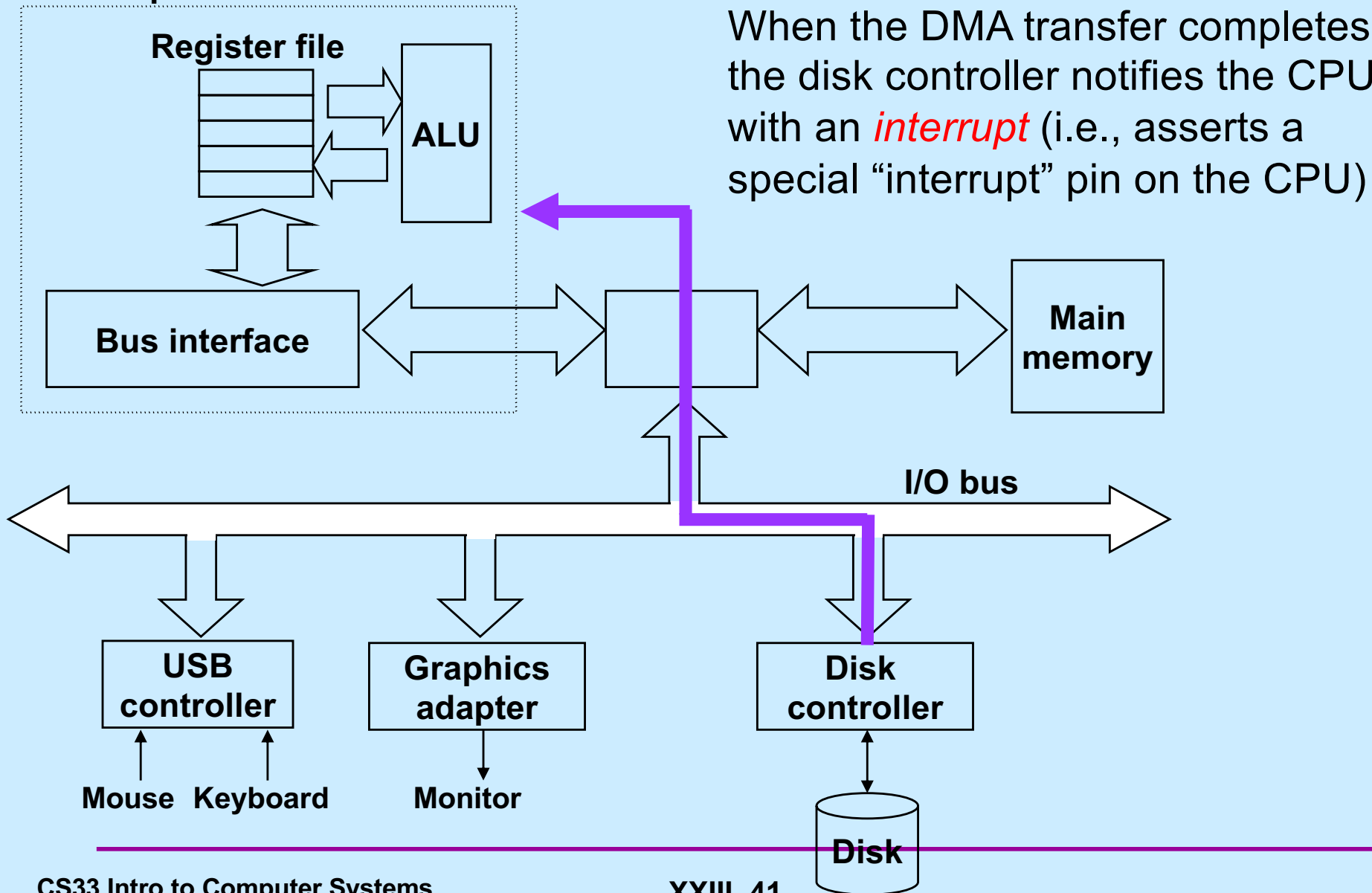




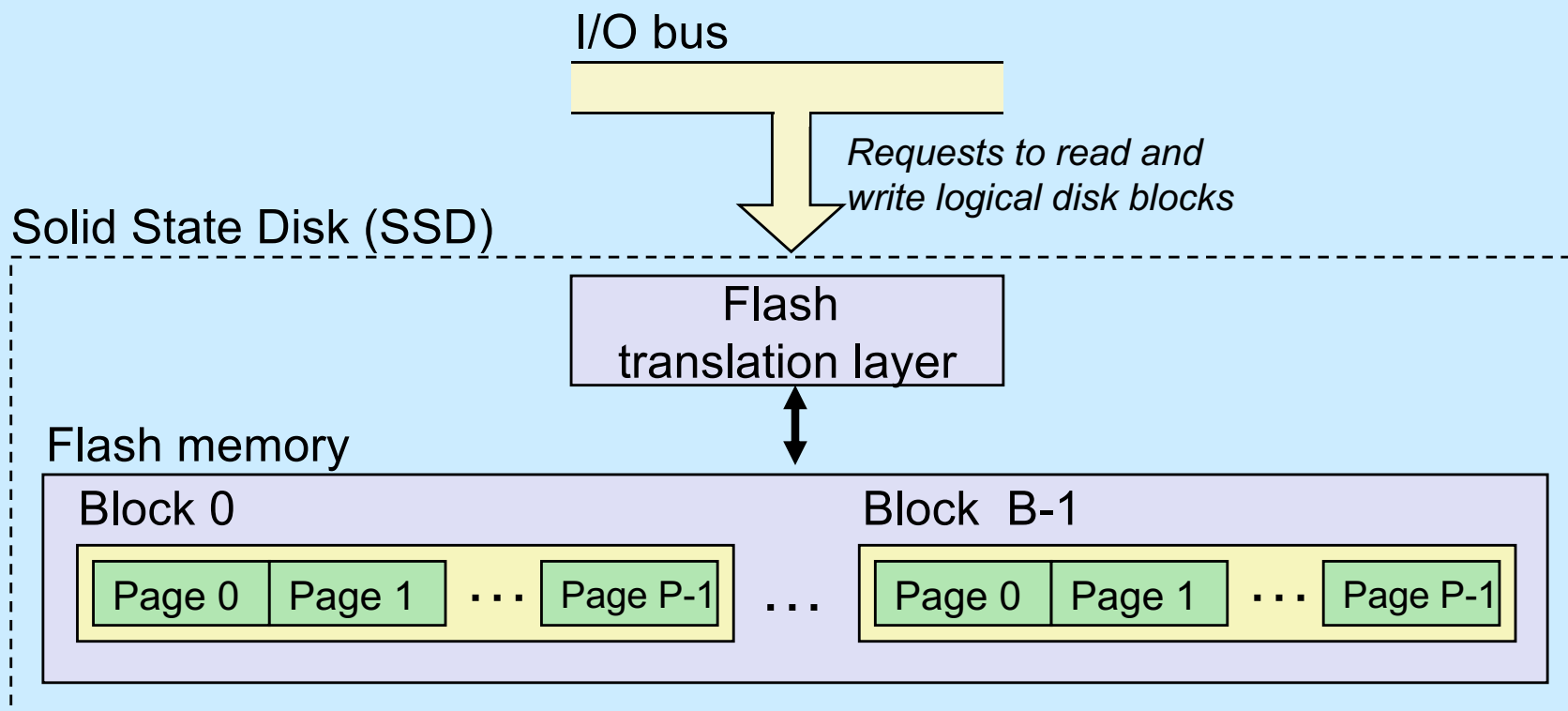
# Reading a Disk Sector (3)

CPU chip

When the DMA transfer completes, the disk controller notifies the CPU with an *interrupt* (i.e., asserts a special “interrupt” pin on the CPU)



# 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

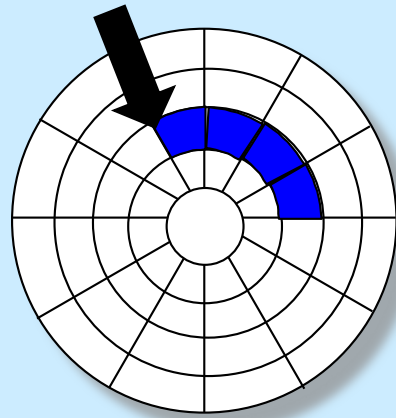
- in 2023, about 1-1.5 times more expensive per byte

- **Applications**

- smart phones, laptops, desktops

# 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**

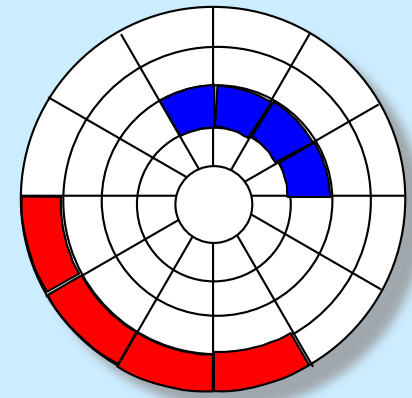


# Quiz 2

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 (within a factor of 2)
- c) much more time



# Quiz 3

**We have two files on the same solid-state disk. Each file's data resides in consecutive blocks. 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 block of the first, then a block of the second, then the first, then the second, etc. Compared to reading them sequentially, this will take**

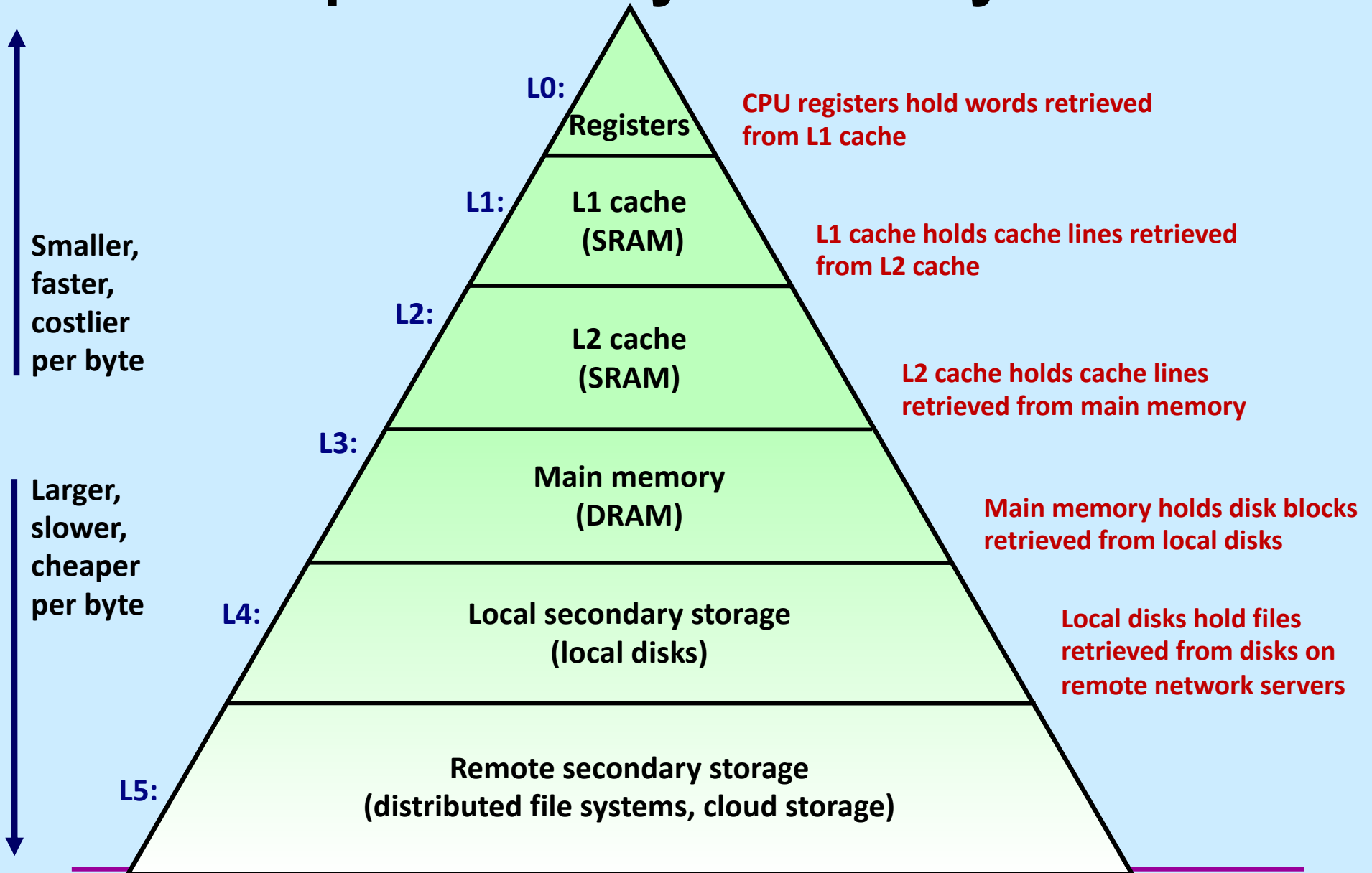
- a) less time**
- b) about the same amount of time  
(within a factor of 2)**
- c) much more time**

# 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 (DRAM) 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

# Disks Are Still Important

- **Cheap**
  - cost/byte less than SSDs
- **(fairly) Reliable**
  - data written to a disk is likely to be there next year
- **Sometimes fast**
  - data in consecutive sectors on a track can be read quickly
- **Sometimes slow**
  - data in randomly scattered sectors takes a long time to read

# Abstraction to the Rescue

- **Programs don't deal with sectors, tracks, and cylinders**
- **Programs deal with *files***
  - **maze.c rather than an ordered collection of sectors**
  - **OS provides the implementation**

# Implementation Problems

- **Speed**
  - **use the hierarchy**
    - » **copy files into RAM, copy back when done**
  - **optimize layout**
    - » **put sectors of a file in consecutive locations**
  - **use parallelism**
    - » **spread file over multiple disks**
    - » **read multiple sectors at once**

# Implementation Problems

- **Reliability**
  - **computer crashes**
    - » what you thought was safely written to the file never made it to the disk — it's still in RAM, which is lost
    - » worse yet, some parts made it back to disk, some didn't
      - you don't know which is which
      - on-disk data structures might be totally trashed
  - **disk crashes**
    - » you had backed it up ... yesterday
  - **you screw up**
    - » you accidentally delete the entire directory containing your shell 1 implementation

# Implementation Problems

- **Reliability solutions**
  - **computer crashes**
    - » **transaction-oriented file systems**
    - » **on-disk data structures always in well defined states**
  - **disk crashes**
    - » **files stored redundantly on multiple disks**
  - **you screw up**
    - » **file system automatically keeps "snapshots" of previous versions of files**