

Selecting a Disk-Scheduling Algorithm

- ▶ SSTF is common and has a natural appeal
- ▶ SCAN and C-SCAN perform better for systems that place a heavy load on the disk.
- ▶ Performance depends on the number and types of requests.
- ▶ Requests for disk service can be influenced by the file-allocation method.
- ▶ The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary.
- ▶ Either SSTF or LOOK is a reasonable choice for the default algorithm.

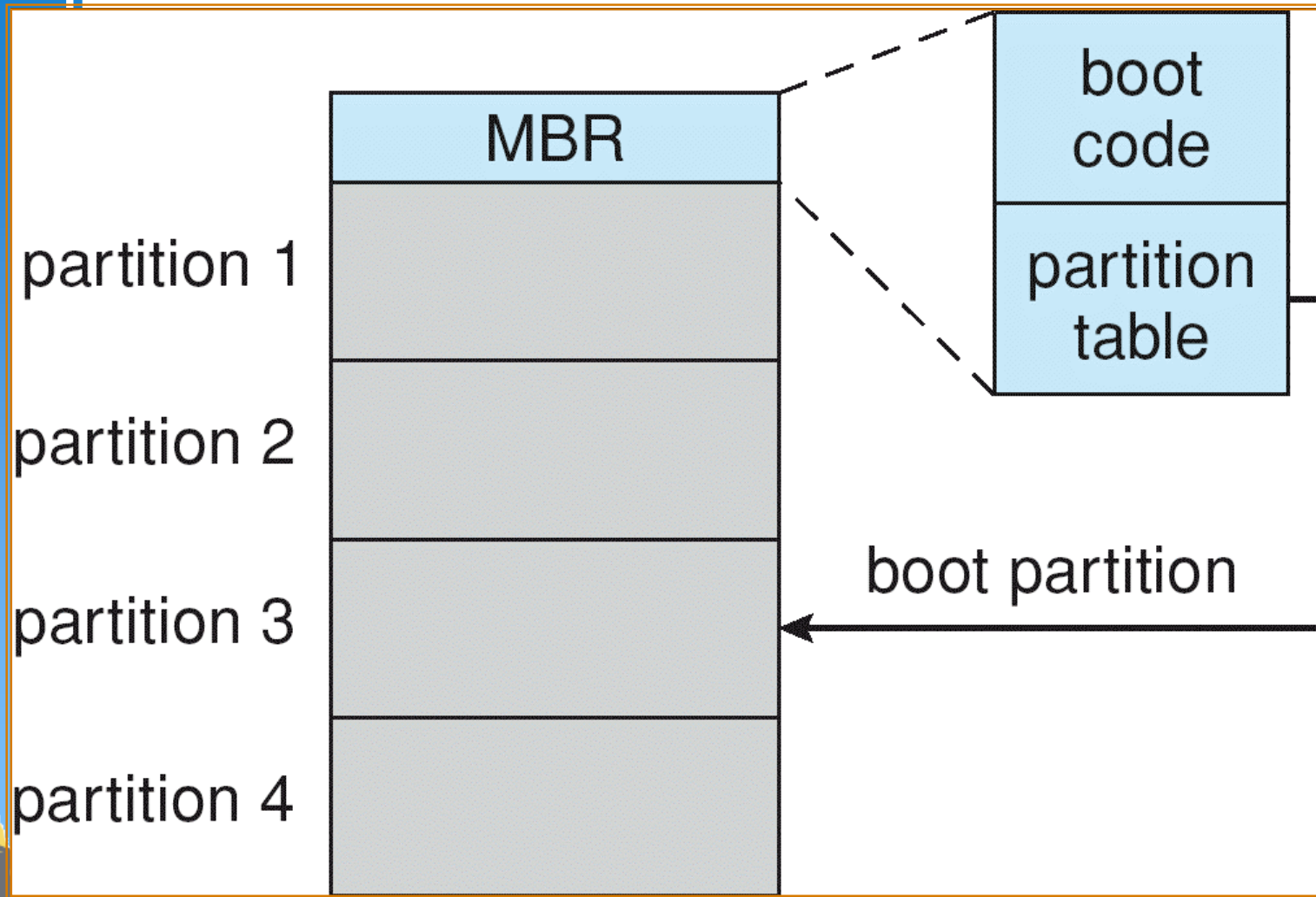


Disk Management

- ▶ *Low-level formatting, or physical formatting* — Dividing a disk into sectors that the disk controller can read and write.
- ▶ To use a disk to hold files, the operating system still needs to record its own data structures on the disk.
 - *Partition* the disk into one or more groups of cylinders.
 - *Logical formatting* or “making a file system”.
- ▶ **Boot block initializes system.**
 - The bootstrap is stored in ROM.
 - *Bootstrap loader* program.
- ▶ Methods such as *sector sparing* used to handle bad blocks.



Booting from a Disk in Windows 2000



Swap-Space Management

- ▶ Swap-space — Virtual memory uses disk space as an extension of main memory.
- ▶ Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition.
- ▶ Swap-space management
 - 4.3BSD allocates swap space when process starts; holds *text segment* (the program) and *data segment*.
 - Kernel uses *swap maps* to track swap-space use.
 - Solaris 2 allocates swap space only when a page is forced out of physical memory, not when the virtual memory page is first created.



RAID Structure

- ▶ **RAID** – multiple disk drives provides **reliability** via **redundancy**.
- ▶ RAID is arranged into six different levels.

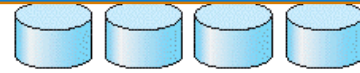


RAID (cont)

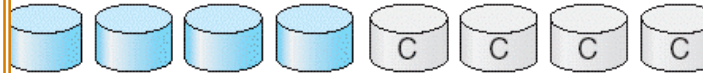
- ▶ Several improvements in disk-use techniques involve the use of multiple disks working cooperatively.
- ▶ Disk striping uses a group of disks as one storage unit.
- ▶ RAID schemes improve performance and improve the reliability of the storage system by storing redundant data.
 - *Mirroring or shadowing* keeps duplicate of each disk.
 - *Block interleaved parity* uses much less redundancy.



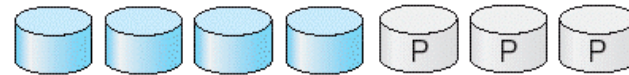
RAID Levels



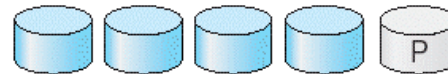
(a) RAID 0: non-redundant striping.



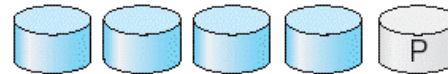
(b) RAID 1: mirrored disks.



(c) RAID 2: memory-style error-correcting codes.



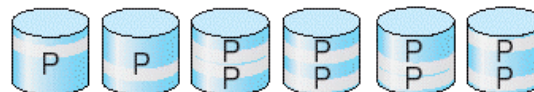
(d) RAID 3: bit-interleaved parity.



(e) RAID 4: block-interleaved parity.



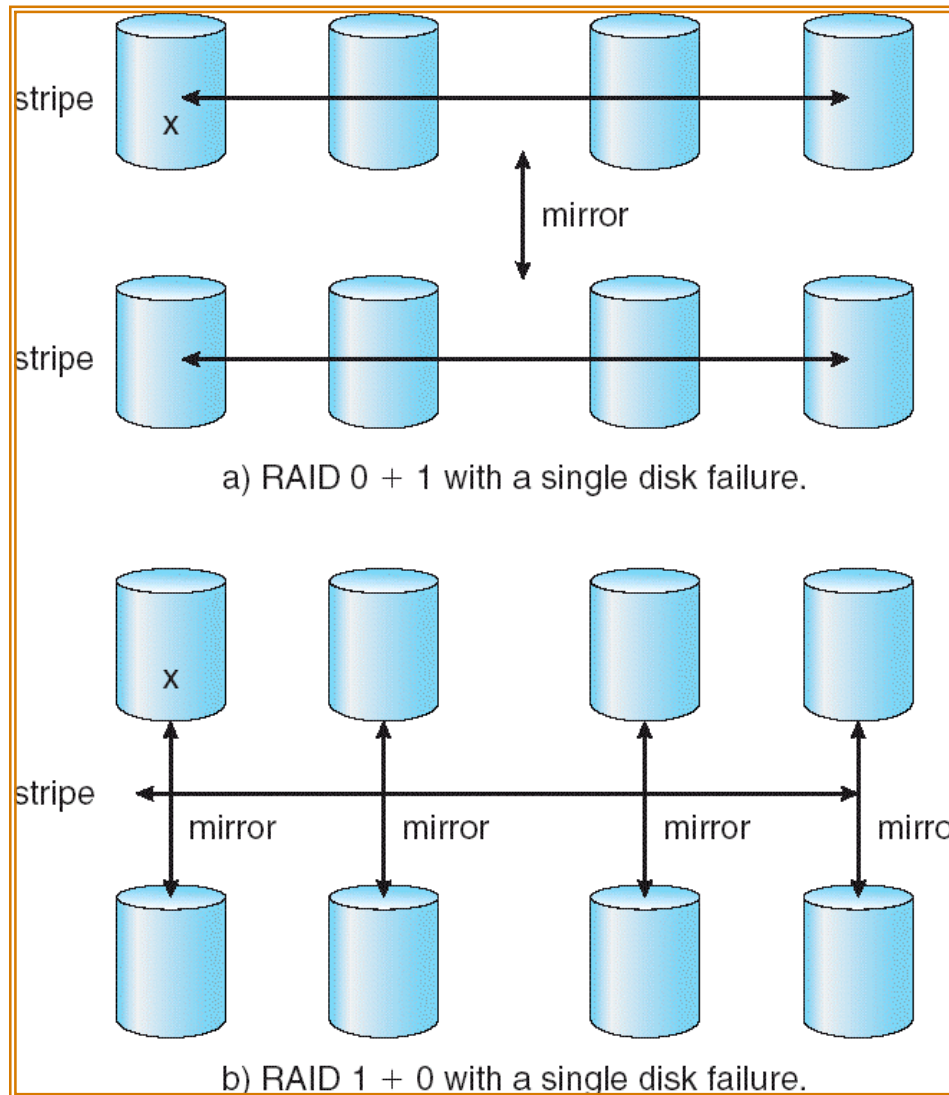
(f) RAID 5: block-interleaved distributed parity.



(g) RAID 6: P + Q redundancy.



RAID (0 + 1) and (1 + 0)



Stable-Storage Implementation

- ▶ Write-ahead log scheme requires stable storage.
- ▶ To implement stable storage:
 - Replicate information on more than one nonvolatile storage media with independent failure modes.
 - Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery.



Chapter 13: I/O Systems - Objectives

- ▶ Explore the structure of an operating system's I/O subsystem
- ▶ Discuss the principles of I/O hardware and its complexity
- ▶ Provide details of the performance aspects of I/O hardware and software

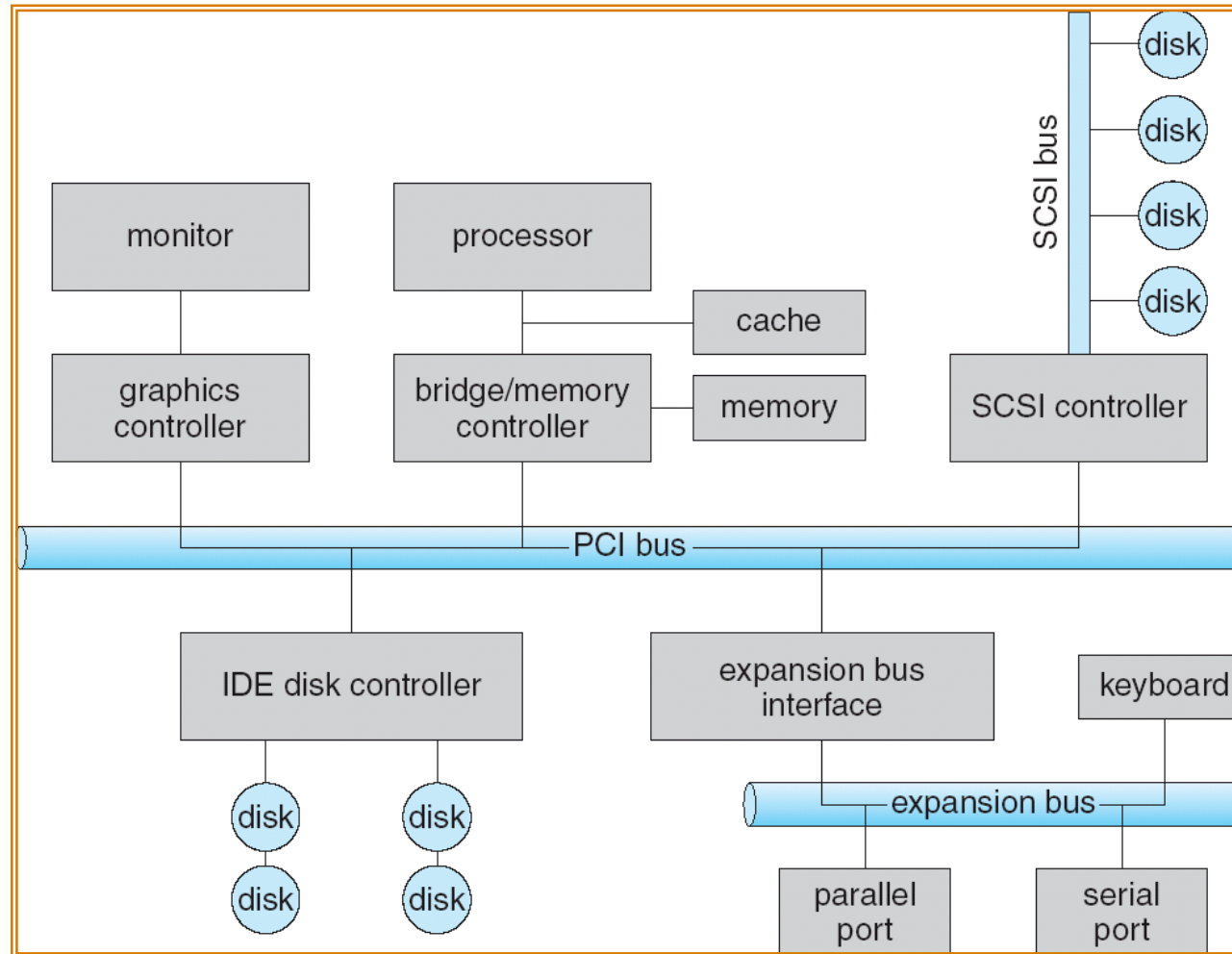


I/O Hardware

- ▶ Incredible variety of I/O devices
- ▶ Common concepts
 - **Port**
 - **Bus (daisy chain or shared direct access)**
 - **Controller (host adapter)**
- ▶ I/O instructions control devices
- ▶ Devices have addresses, used by
 - Direct I/O instructions
 - **Memory-mapped I/O**



A Typical PC Bus Structure



Polling

- ▶ Determines state of device
 - command-ready
 - busy
 - Error
- ▶ **Busy-wait** cycle to wait for I/O from device

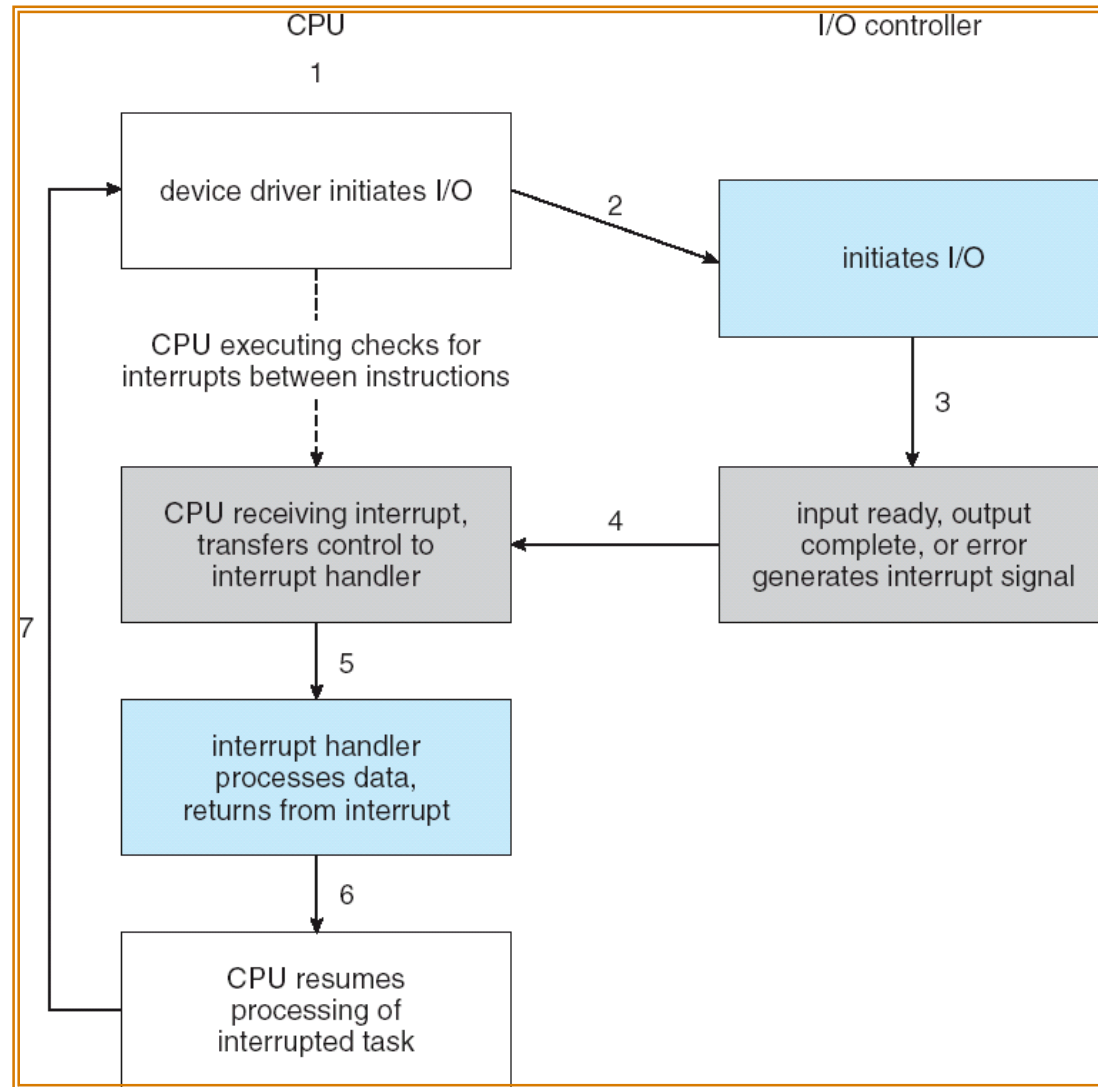


Interrupts

- ▶ **CPU Interrupt-request line** triggered by I/O device
- ▶ **Interrupt handler** receives interrupts
- ▶ **Maskable** to ignore or delay some interrupts
- ▶ Interrupt vector to dispatch interrupt to correct handler
 - Based on priority
 - Some **nonmaskable**
- ▶ Interrupt mechanism also used for exceptions



Interrupt-Driven I/O Cycle

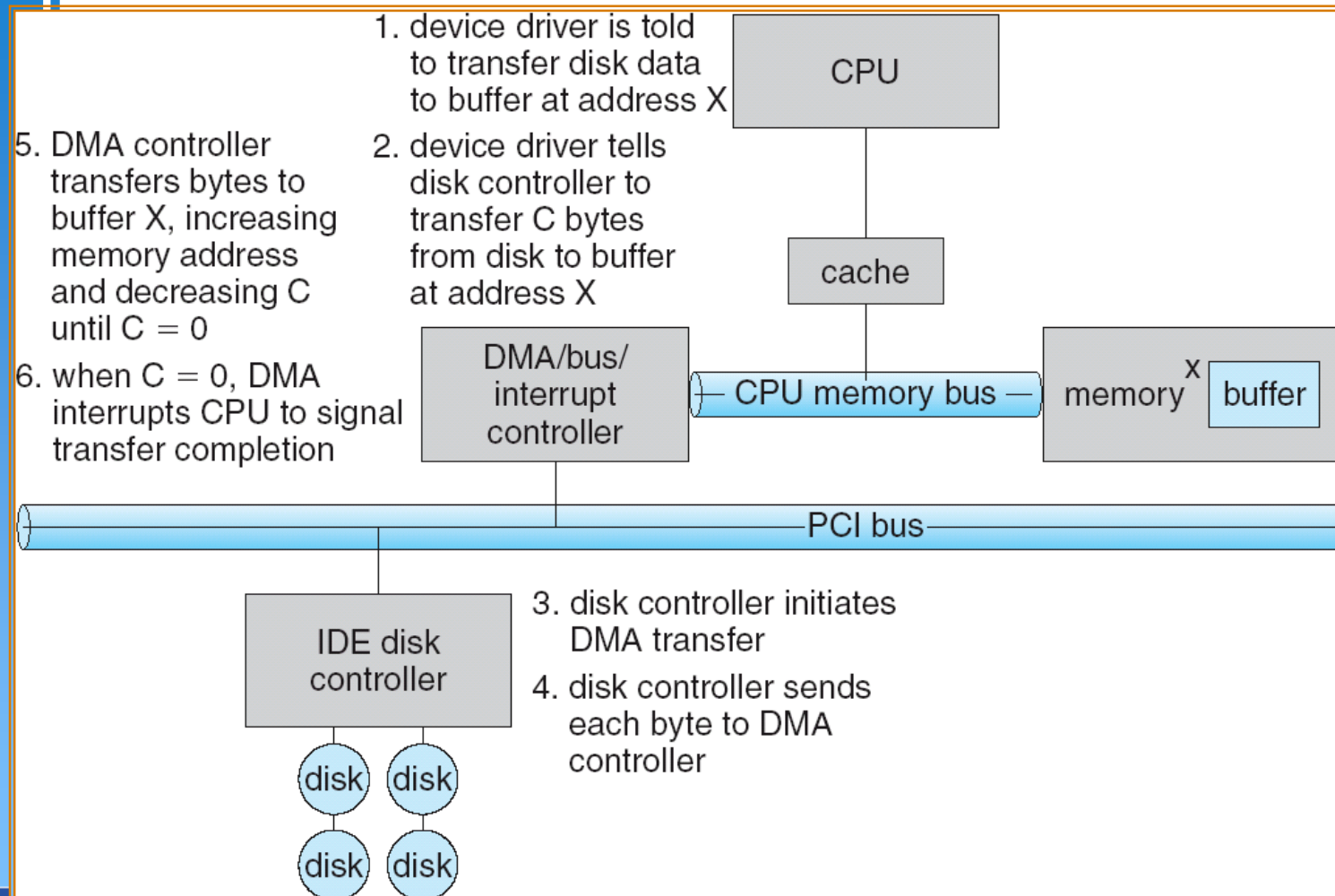


Direct Memory Access

- ▶ Used to avoid **programmed I/O** for large data movement
- ▶ Requires **DMA** controller
- ▶ Bypasses CPU to transfer data directly between I/O device and memory



Six Step Process to Perform DMA Transfer

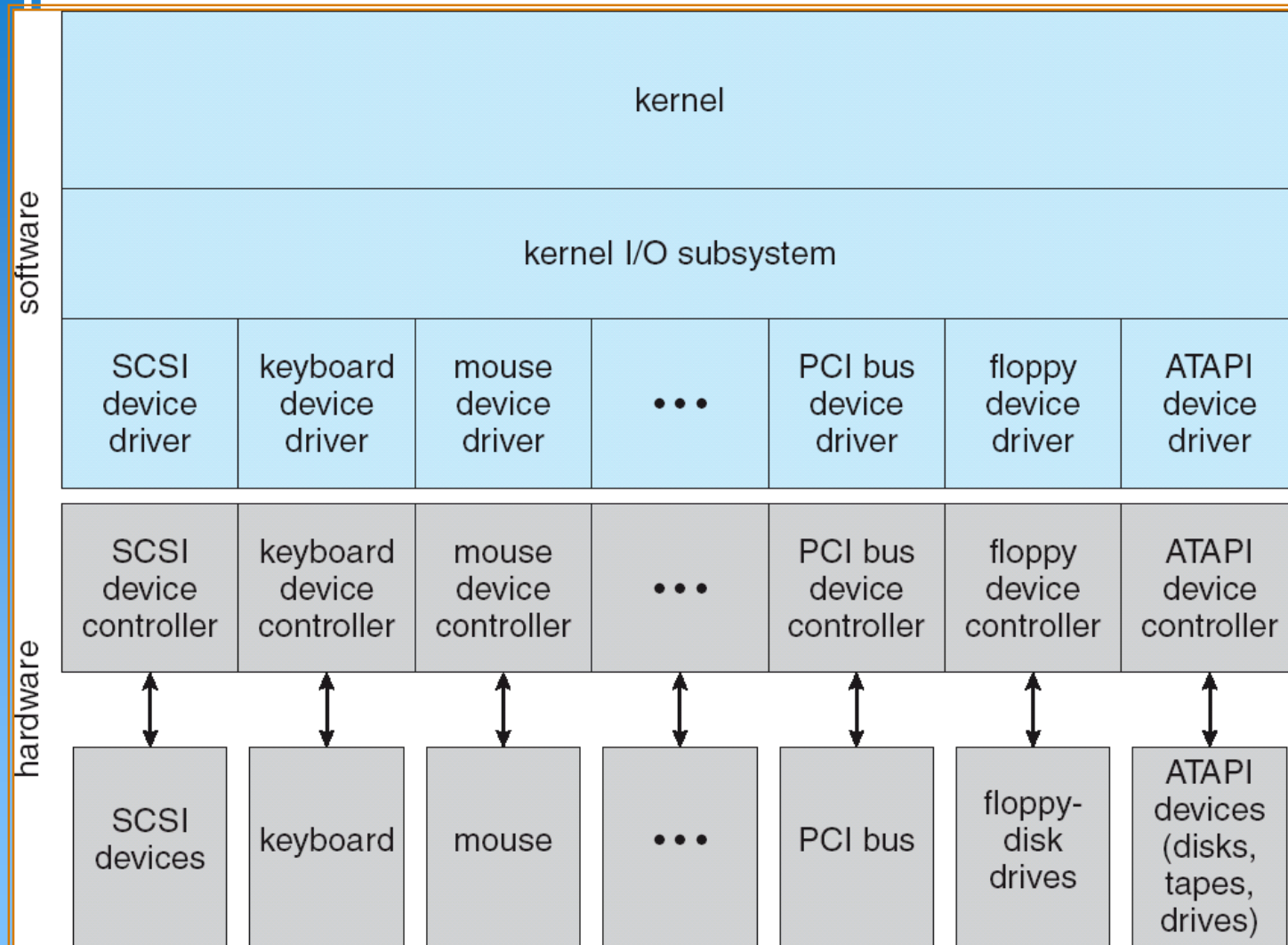


Application I/O Interface

- ▶ I/O system calls encapsulate device behaviors in generic classes
- ▶ Device-driver layer hides differences among I/O controllers from kernel
- ▶ Devices vary in many dimensions
 - **Character-stream or block**
 - **Sequential or random-access**
 - **Sharable or dedicated**
 - **Speed of operation**
 - **read-write, read only, or write only**



A Kernel I/O Structure



Characteristics of I/O Devices

aspect	variation	example
data-transfer mode	character block	terminal disk
access method	sequential random	modem CD-ROM
transfer schedule	synchronous asynchronous	tape keyboard
sharing	dedicated sharable	tape keyboard
device speed	latency seek time transfer rate delay between operations	
I/O direction	read only write only read–write	CD-ROM graphics controller disk



Block and Character Devices

- ▶ Block devices include disk drives
 - Commands include read, write, seek
 - Raw I/O or file-system access
 - Memory-mapped file access possible
- ▶ Character devices include keyboards, mice, serial ports
 - Commands include `get`, `put`
 - Libraries layered on top allow line editing

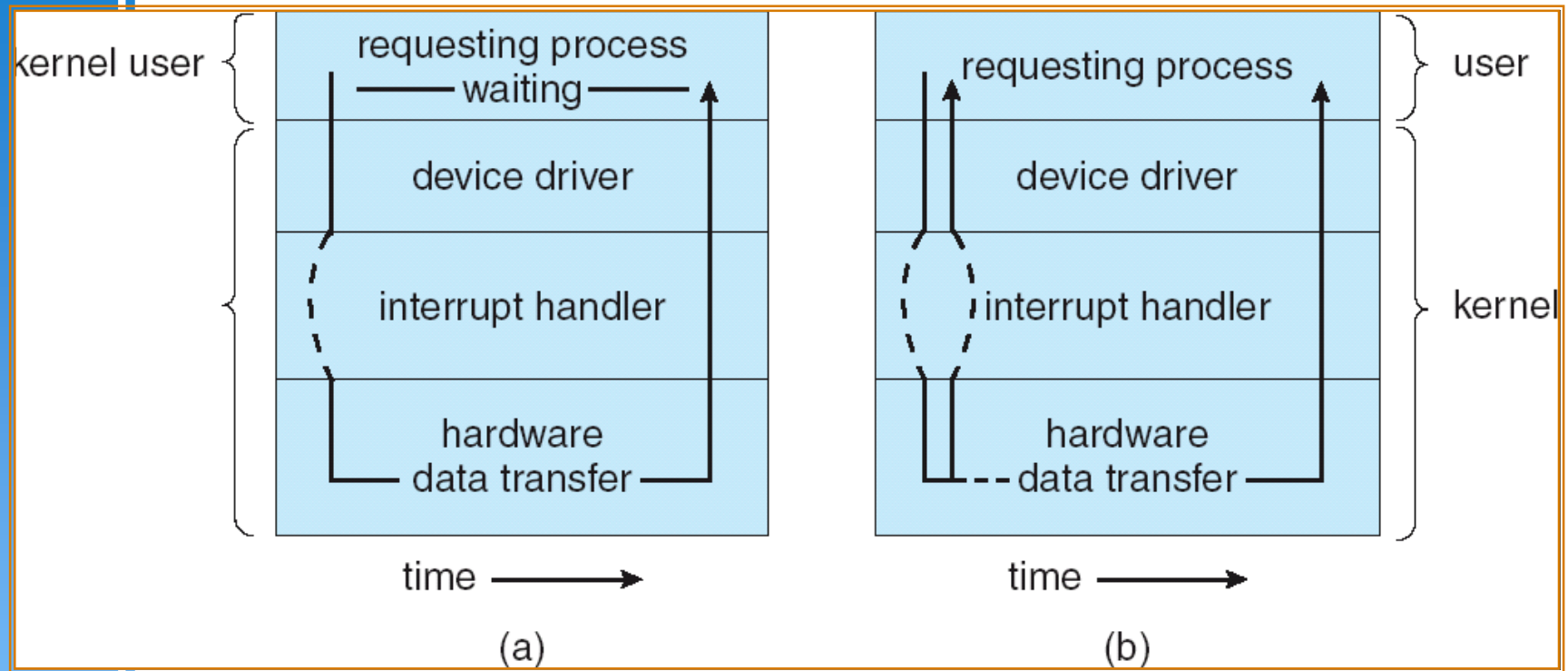


Blocking and Nonblocking I/O

- ▶ **Blocking** - process suspended until I/O completed
 - Easy to use and understand
 - Insufficient for some needs
- ▶ **Nonblocking** - I/O call returns as much as available
 - User interface, data copy (buffered I/O)
 - Implemented via multi-threading
 - Returns quickly with count of bytes read or written
- ▶ **Asynchronous** - process runs while I/O executes
 - Difficult to use
 - I/O subsystem signals process when I/O completed



Two I/O Methods



Synchronous

Asynchronous



Kernel I/O Subsystem

- ▶ Scheduling
 - Some I/O request ordering via per-device queue
 - Some OSs try fairness
- ▶ Buffering - store data in memory while transferring between devices
 - To cope with device speed mismatch
 - To cope with device transfer size mismatch
 - To maintain “copy semantics”



Kernel I/O Subsystem

- ▶ **Caching** - fast memory holding copy of data
 - Always just a copy
 - Key to performance
- ▶ **Spooling** - hold output for a device
 - If device can serve only one request at a time
 - i.e., Printing
- ▶ **Device reservation** - provides exclusive access to a device
 - System calls for allocation and deallocation
 - Watch out for deadlock



Error Handling

- ▶ OS can recover from disk read, device unavailable, transient write failures
- ▶ Most return an error number or code when I/O request fails
- ▶ System error logs hold problem reports



I/O Protection

- ▶ User process may accidentally or purposefully attempt to disrupt normal operation via illegal I/O instructions
 - All I/O instructions defined to be privileged
 - I/O must be performed via system calls
 - Memory-mapped and I/O port memory locations must be protected too



Kernel Data Structures

- ▶ Kernel keeps state info for I/O components, including open file tables, network connections, character device state
- ▶ Many, many complex data structures to track buffers, memory allocation, “dirty” blocks
- ▶ Some use object-oriented methods and message passing to implement I/O



I/O Requests to Hardware Operations

- ▶ Consider reading a file from disk for a process:
 - Determine device holding file
 - Translate name to device representation
 - Physically read data from disk into buffer
 - Make data available to requesting process
 - Return control to process



Improving Performance

- ▶ Reduce number of context switches
- ▶ Reduce data copying
- ▶ Reduce interrupts by using large transfers, smart controllers, polling
- ▶ Use DMA
- ▶ Balance CPU, memory, bus, and I/O performance for highest throughput

