# 13. I/O Systems

*Sungyoung Lee*

*College of Engineering*

*KyungHee University*
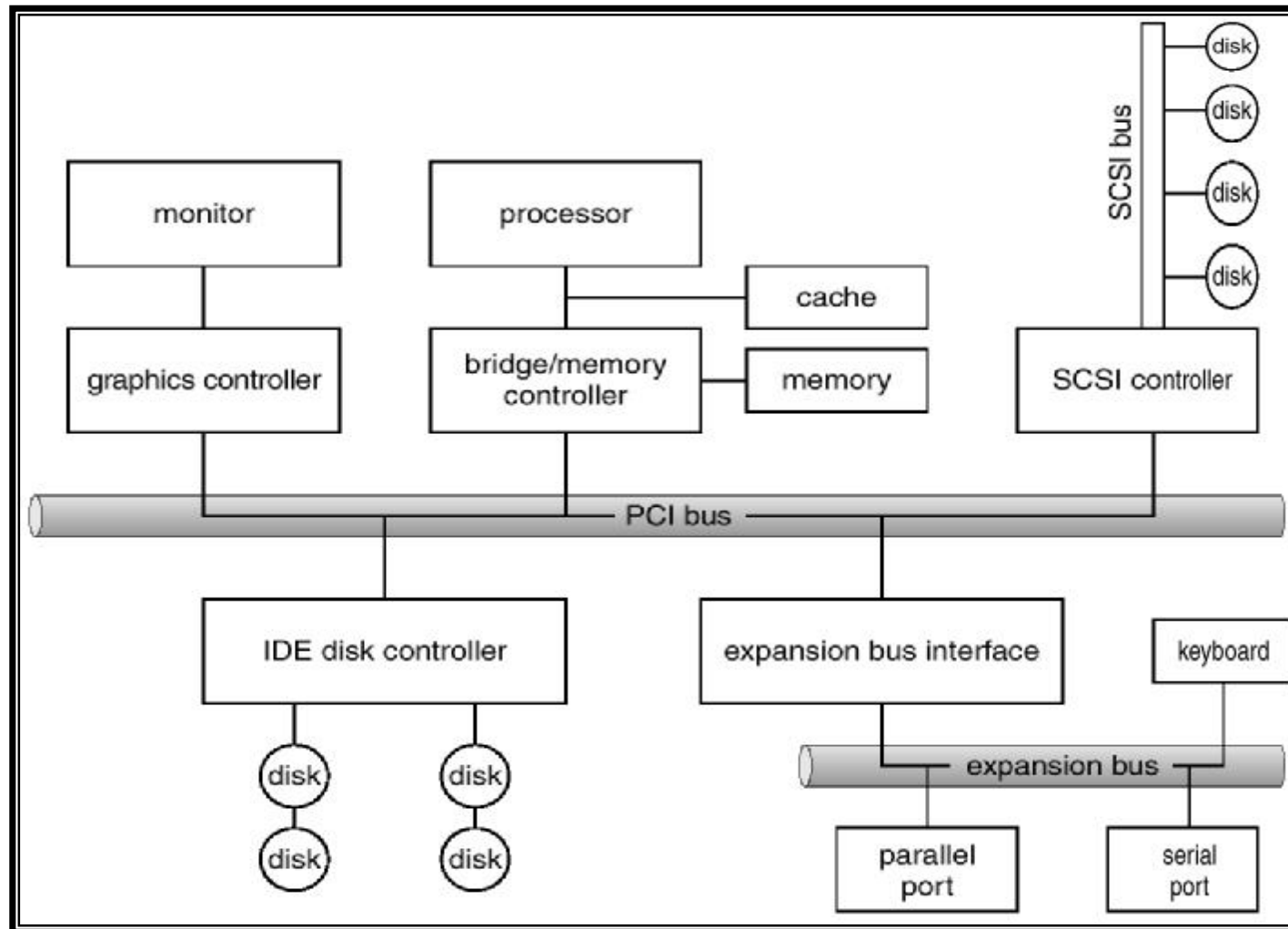
# Contents

# I/O Hardware

n Incredible variety of I/O devices

n Common concepts
   ü Port
   ü Bus (daisy chain or shared direct access)
   ü Controller (host adapter)

n I/O instructions control devices

n Devices have addresses, used by
   ü Direct I/O instructions
   ü Memory-mapped I/O

**n** Device controller (or host adapter)

- ü I/O devices have components:
  - § Mechanical component
  - § Electronic component
- ü The electronic component is the device controller
  - § May be able to handle multiple devices
- ü Controller's tasks
  - § Convert serial bit stream to block of bytes
  - § Perform error correction as necessary
  - § Make available to main memory

**n** Use special I/O instructions to an I/O port address

| I/O address range (hexadecimal) | device |
| --- | --- |
| 000-00F | DMA controller |
| 020-021 | interrupt controller |
| 040-043 | timer |
| 200-20F | game controller |
| 2F8-2FF | serial port (secondary) |
| 320-32F | hard-disk controller |
| 378-37F | parallel port |
| 3D0-3DF | graphics controller |
| 3F0-3F7 | diskette-drive controller |
| 3F8-3FF | serial port (primary) |

- **n** The device control registers are mapped into the address space of the processor
  - ü The CPU executes I/O requests using the standard data transfer instructions
- **n** I/O device drivers can be written entirely in C
- **n** No special protection mechanism is needed to keep user processes from performing I/O
  - ü Can give a user control over specific devices but not others by simply including the desired pages in its page table
- **n** Reading a device register and testing its value is done with a single instruction
- **n** Memory-mapped regions should be uncacheable
- **n** Memory-mapped device register is vulnerable to accidental modification through the use of incorrect pointers
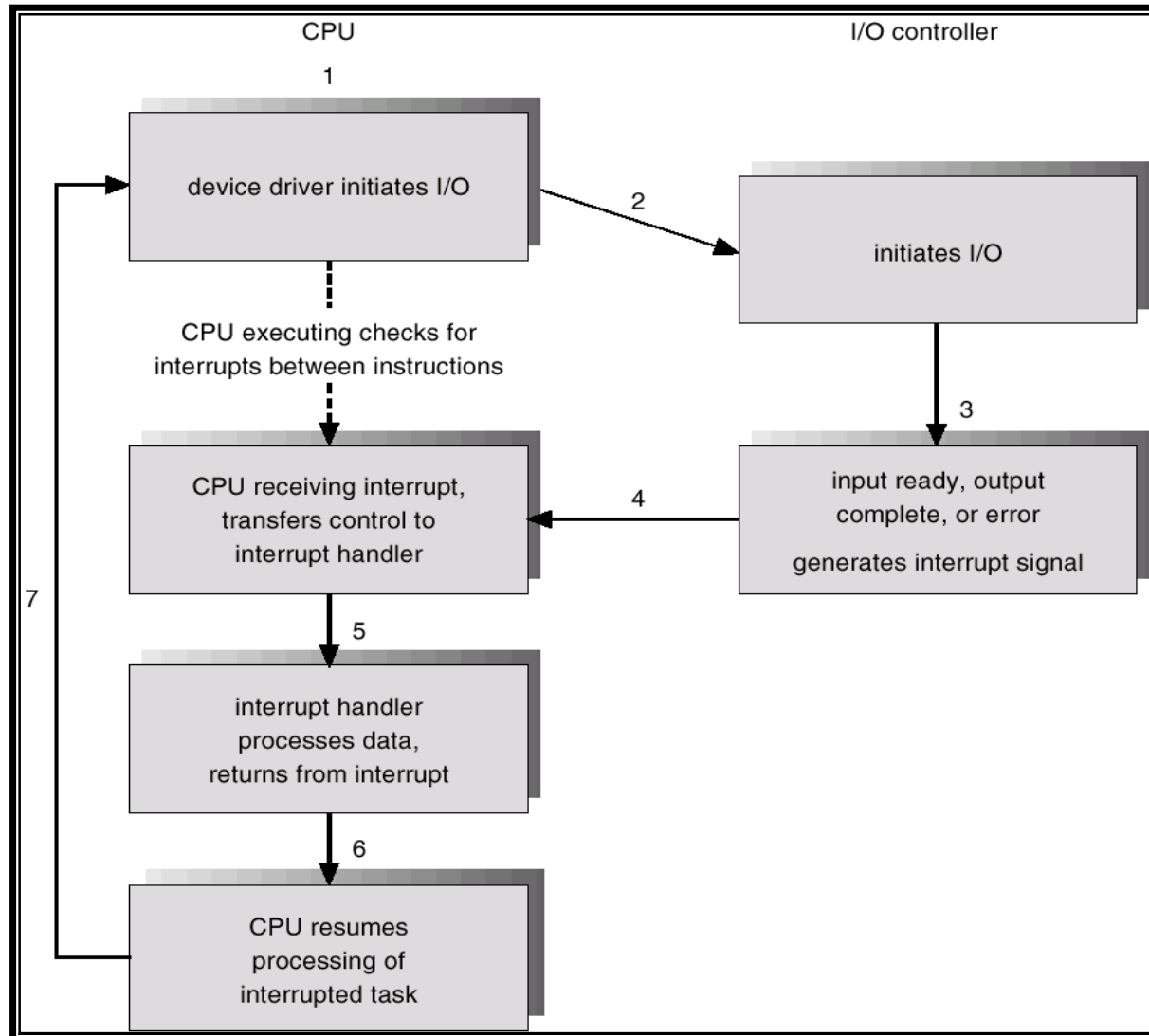  - ü Protected memory helps to reduce this risk

# Polling

**n** Determines state of device

  **ü** command-ready

  **ü** busy

  **ü** Error

**n** Busy-wait cycle to wait for I/O from device

# Interrupts

- **n** CPU Interrupt request line triggered by I/O device

- **n** Interrupt handler receives interrupts

- **n** Maskable to ignore or delay some interrupts

- **n** Interrupt vector to dispatch interrupt to correct handler
  - ü Based on priority
  - ü Some unmaskable

- **n** Interrupt mechanism also used for exceptions

# Interrupt-Driven I/O Cycle

# Intel Pentium Processor Event-Vector Table

| vector number | description |
| --- | --- |
| 0 | divide error |
| 1 | debug exception |
| 2 | null interrupt |
| 3 | breakpoint |
| 4 | INTO-detected overflow |
| 5 | bound range exception |
| 6 | invalid opcode |
| 7 | device not available |
| 8 | double fault |
| 9 | coprocessor segment overrun (reserved) |
| 10 | invalid task state segment |
| 11 | segment not present |
| 12 | stack fault |
| 13 | general protection |
| 14 | page fault |
| 15 | (Intel reserved, do not use) |
| 16 | floating-point error |
| 17 | alignment check |
| 18 | machine check |
| 19Ð31 | (Intel reserved, do not use) |
| 32Ð255 | maskable interrupts |

**n** Polled I/O

   ü CPU asks ("polls") devices if need attention

       § ready to receive a command

       § command status, etc.

   ü Advantages

       § Simple

       § Software is in control

       § Efficient if CPU finds a device to be ready soon

   ü Disadvantages

       § Inefficient in non-trivial system (high CPU utilization)

       § Low priority devices may never be serviced
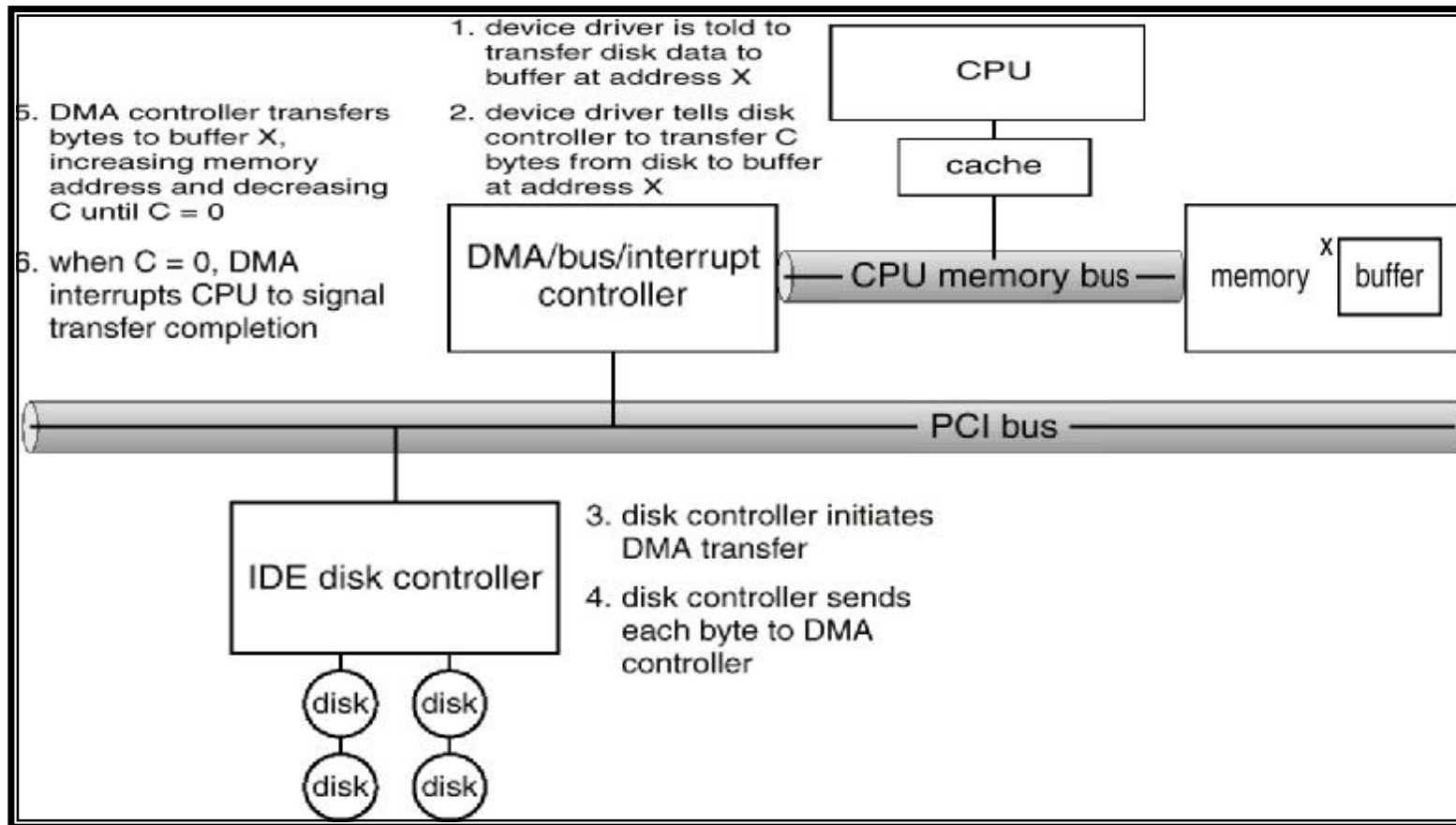
# Polling vs. Interrupts (Cont'd)

**n** Interrupt-driven I/O
- ü I/O devices request interrupt when need attention
- ü Interrupt service routines specific to each device are invoked
- ü Interrupts can be shared between multiple devices
- ü Advantages
  - § CPU only attends to device when necessary
  - § More efficient than polling in general
- ü Disadvantages
  - § Excess interrupts slow (or prevent) program execution
  - § Overheads (may need 1 interrupt per byte transferred)

# Direct Memory Access

**n** Used to avoid programmed I/O for large data movement

    **ü** Programmed I/O?

**n** Requires DMA controller

**n** Bypasses CPU to transfer data directly between I/O device and memory

**n** DMA modes

(1) Cycle stealing

§ The DMA controller sneaks in and steals an occasional bus cycle from the CPU once in a while, delaying it slightly

(2) Burst mode

§ The DMA controller acquires the bus, issues a series of transfers, then releases the bus

§ More efficient than cycle stealing: acquiring the bus takes time and multiple words can be transferred for the price of one bus acquisition

§ It can block the CPU and other devices too long

**n  Addressing in DMA**

(1) Physical address

§  OS converts the virtual address of the intended memory buffer into a physical address and writes it into DMA controller's address register

(2) Virtual address

§  The DMA controller must use the MMU to have the virtual-to-physical translation done

§  Not common: only when the MMU is part of the memory rather than part of the CPU

ü  In any case, the target memory region should be pinned (not paged out) during DMA

n **DMA types**

(1) Sequential DMA

- § Data temporarily stored in DMA controller
- § Requires an extra bus cycle per word transferred
- § More flexible in that it can also perform device-to-device copies and even memory-to-memory copies
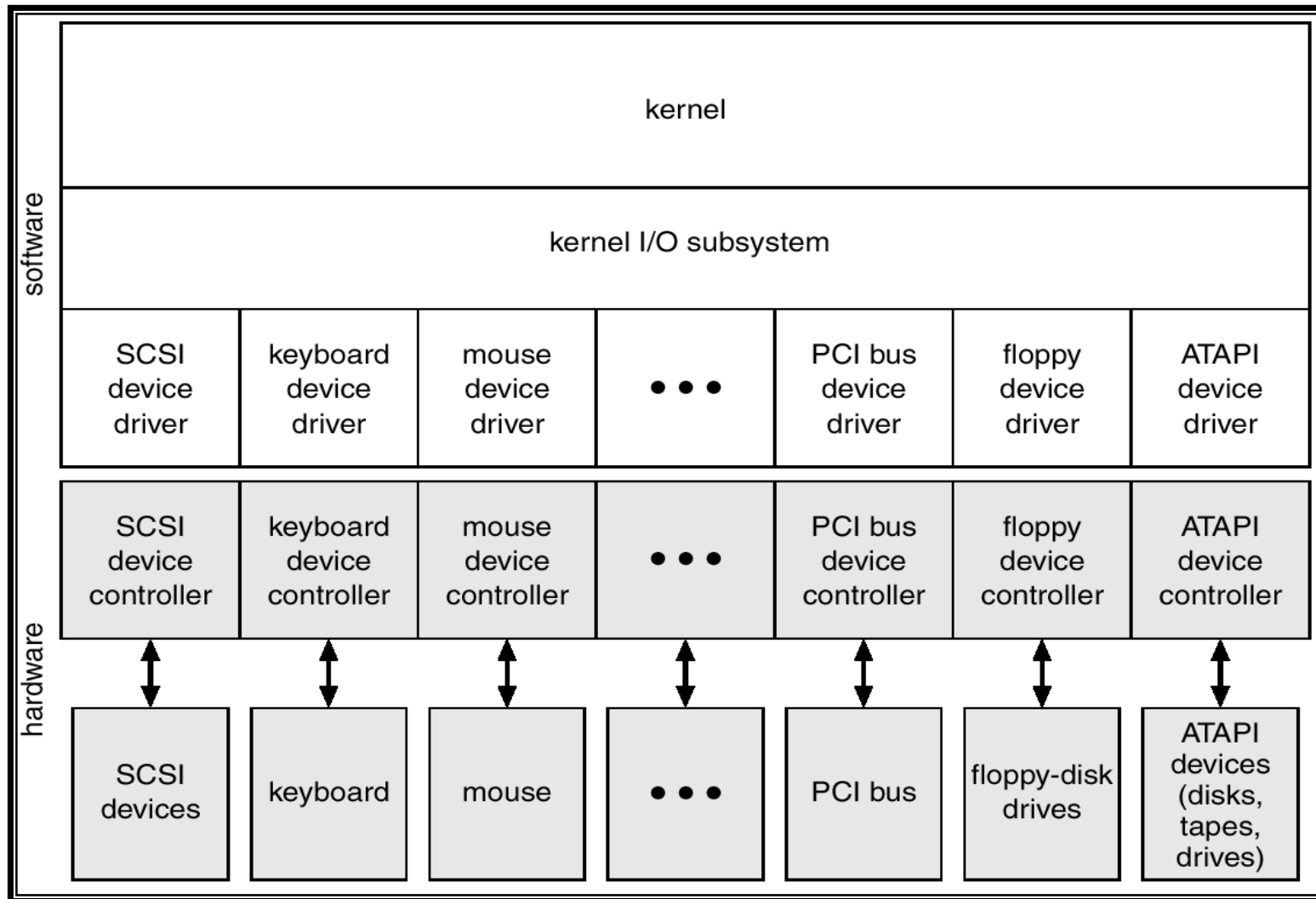
(2) Simultaneous DMA (or fly-by mode)

- § The DMA controller tells the device controller to transfer the data directly to main memory

# Application I/O Interface

**n** I/O system calls encapsulate device behaviors in generic classes

**n** Device-driver layer hides differences among I/O controllers from kernel

**n** Devices vary in many dimensions
  - ü Character-stream or block
  - ü Sequential or random-access
  - ü Sharable or dedicated
  - ü Speed of operation
  - ü read-write, read only, or write only

# A Kernel I/O Structure

# Characteristics of I/O Devices

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

# Block and Character Devices

**n** Block devices include disk drives
- **ü** Commands include `read, write, seek`
- **ü** Raw I/O or file-system access
- **ü** Memory-mapped file access possible

**n** Character devices include keyboards, mice, serial ports
- **ü** Commands include `get, put`
- **ü** Libraries layered on top allow line editing

# Network Devices

**n**  Varying enough from block and character to have own interface

**n**  Unix and Windows NT/9*i*/2000 include socket interface
  - ü Separates network protocol from network operation
  - ü Includes `select` functionality

**n**  Approaches vary widely (pipes, FIFOs, streams, queues, mailboxes)

# Clocks and Timers

**n**  Provide current time, elapsed time, timer

**n**  If programmable interval time used for timings, periodic interrupts

**n**  `ioctl` (on UNIX) covers odd aspects of I/O such as clocks and timers

# Blocking and Nonblocking I/O

**n** Blocking - process suspended until I/O completed
- ü Easy to use and understand
- ü Insufficient for some needs

**n** Nonblocking - I/O call returns as much as available
- ü User interface, data copy (buffered I/O)
- ü Implemented via multi-threading
- ü Returns quickly with count of bytes read or written

**n** Asynchronous - process runs while I/O executes
- ü Difficult to use
- ü I/O subsystem signals process when I/O completed

# Kernel I/O Subsystem

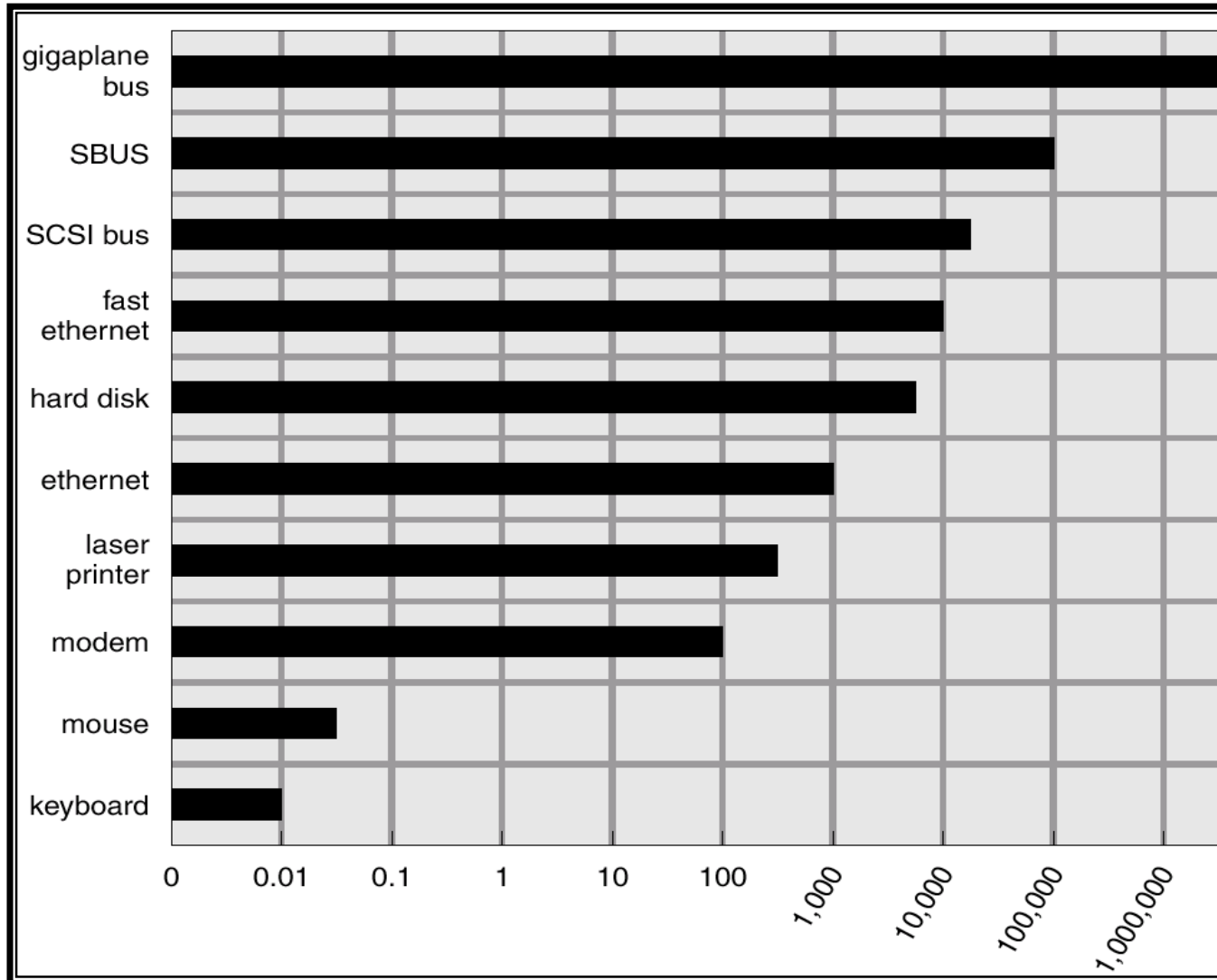**n**  Scheduling

 ü  Some I/O request ordering via per-device queue

 ü  Some OSs try fairness

**n**  Buffering - store data in memory while transferring between devices

 ü  To cope with device speed mismatch

 ü  To cope with device transfer size mismatch

 ü  To maintain "copy semantics"

# Kernel I/O Subsystem

**n** Caching - fast memory holding copy of data
- ü Always just a copy
- ü Key to performance

**n** Spooling - hold output for a device
- ü If device can serve only one request at a time
- ü i.e., Printing

**n** Device reservation - provides exclusive access to a device
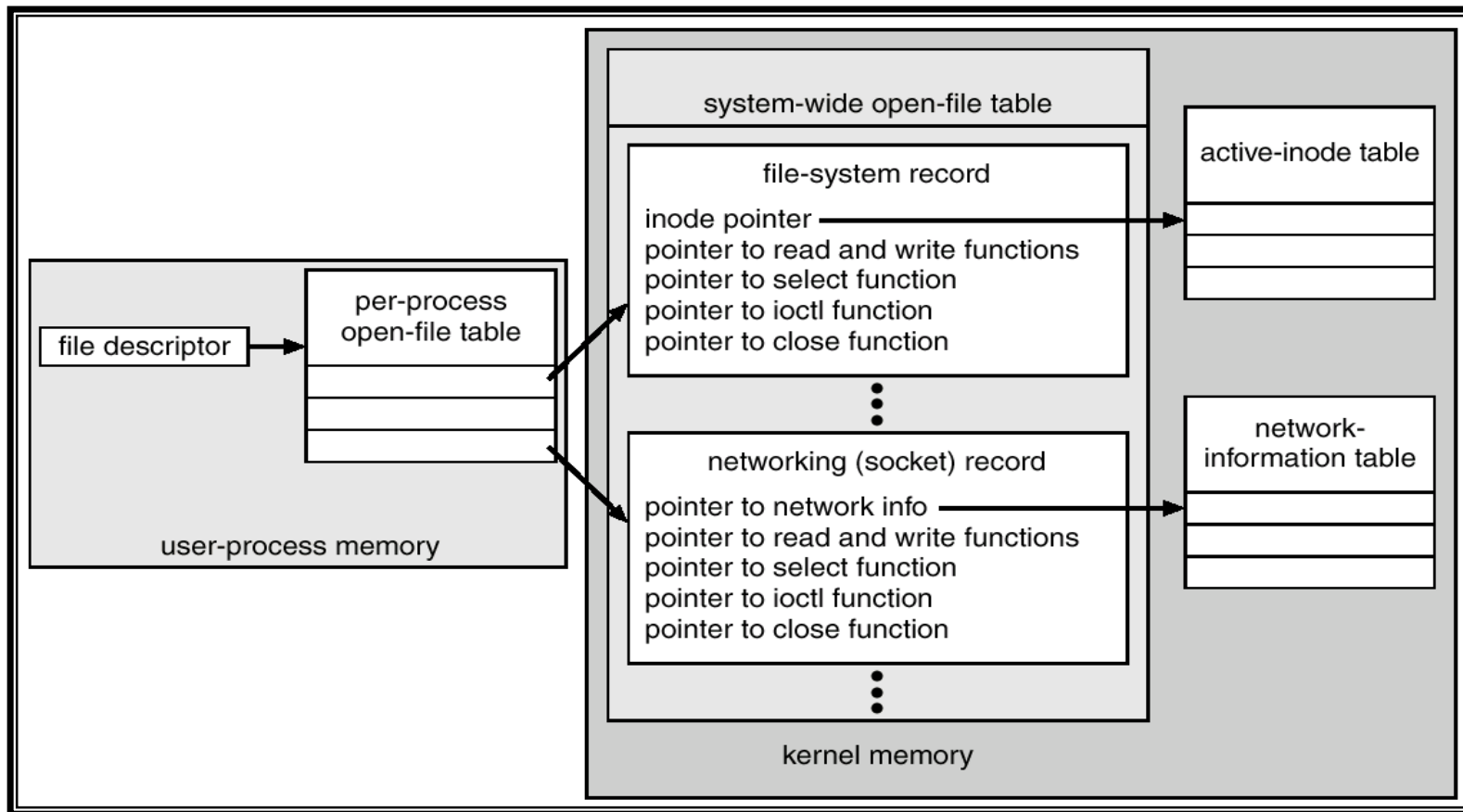- ü System calls for allocation and deallocation
- ü Watch out for deadlock

# Error Handling

- **n** OS can recover from disk read, device unavailable, transient write failures

- **n** Most return an error number or code when I/O request fails

- **n** System error logs hold problem reports

# Kernel Data Structures

n Kernel keeps state info for I/O components, including open file tables, network connections, character device state

n Many, many complex data structures to track buffers, memory allocation, "dirty" blocks

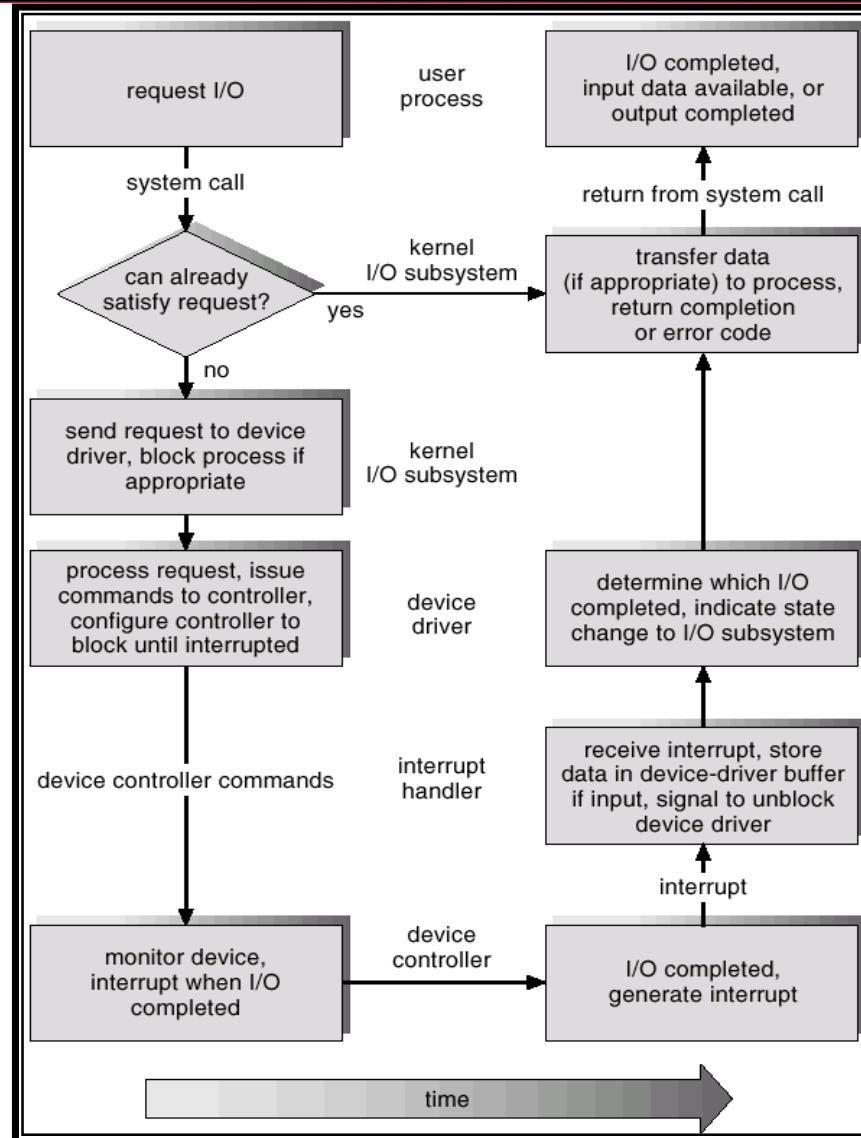n Some use object-oriented methods and message passing to implement I/O

# I/O Requests to Hardware Operations

**n** Consider reading a file from disk for a process:

- ü Determine device holding file
- ü Translate name to device representation
- ü Physically read data from disk into buffer
- ü Make data available to requesting process
- ü Return control to process

# STREAMS

**n** **STREAM**
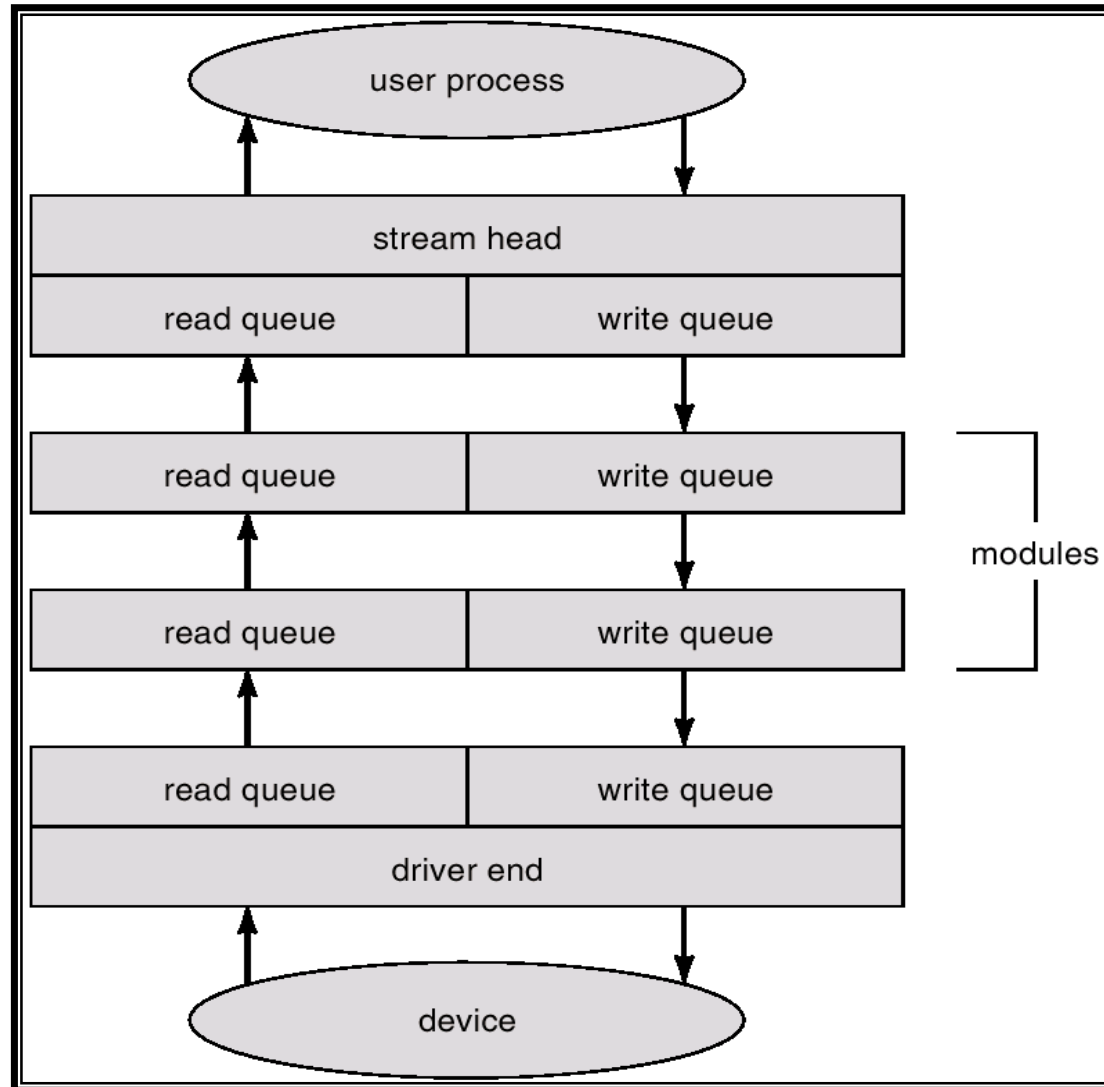  - ü  a full-duplex communication channel between a user-level process and a device

**n** A STREAM consists of:
  - ü  **STREAM head** interfaces with the user process
  - ü  **driver end** interfaces with the device
  - ü  zero or more STREAM modules between them

**n** Each module contains a **read  queue** and a **write queue**

**n** Message passing is used to communicate between queues

# The STREAMS Structure

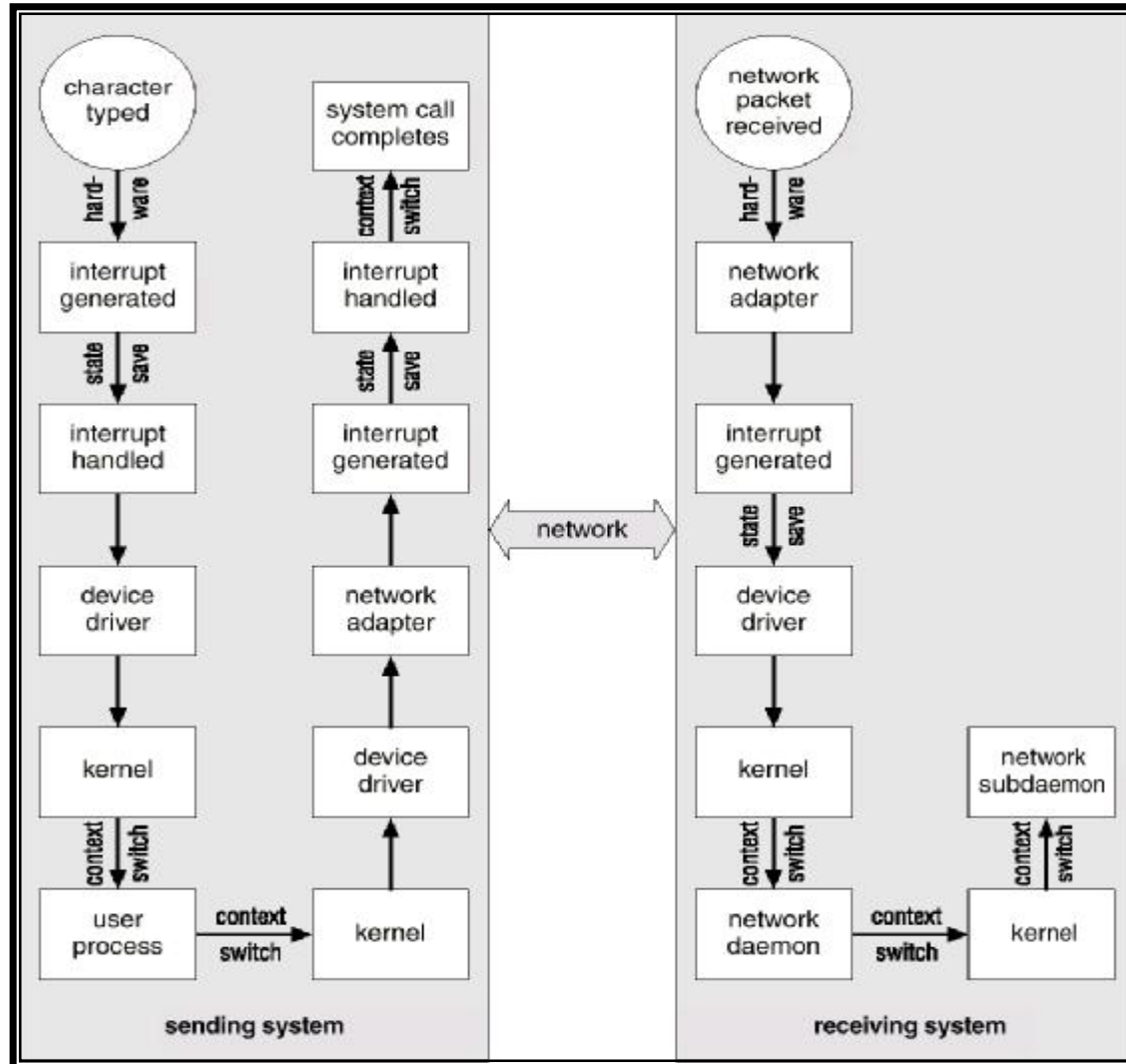# Performance

n  I/O a major factor in system performance:

  ü Demands CPU to execute device driver, kernel I/O code

  ü Context switches due to interrupts

  ü Data copying

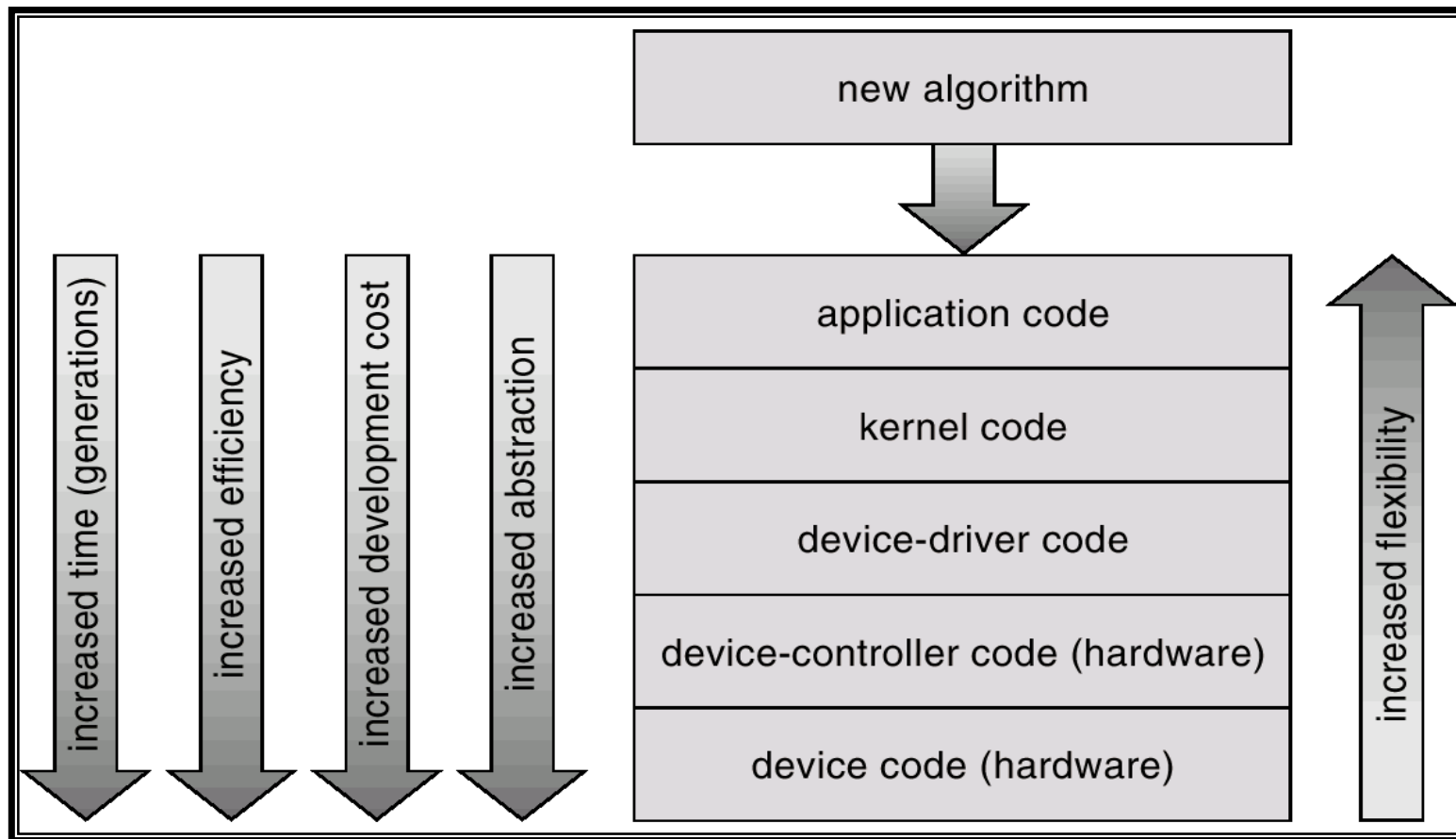  ü Network traffic especially stressful

# Intercomputer Communications

# Improving Performance

**n** Reduce number of context switches

**n** Reduce data copying

**n** Reduce interrupts by using large transfers, smart controllers, polling

**n** Use DMA

**n** Balance CPU, memory, bus, and I/O performance for highest throughput

# Device-Functionality Progression

**n** Goals

- ü Device independence
    - § Programs can access any I/O device without specifying device in advance
- ü Uniform naming
    - § Name of a file or device should simply be a string or an integer
- ü Error handling
    - § Handle as close to the hardware as possible
- ü Synchronous vs. asynchronous
    - § blocked transfers vs. interrupt-driven
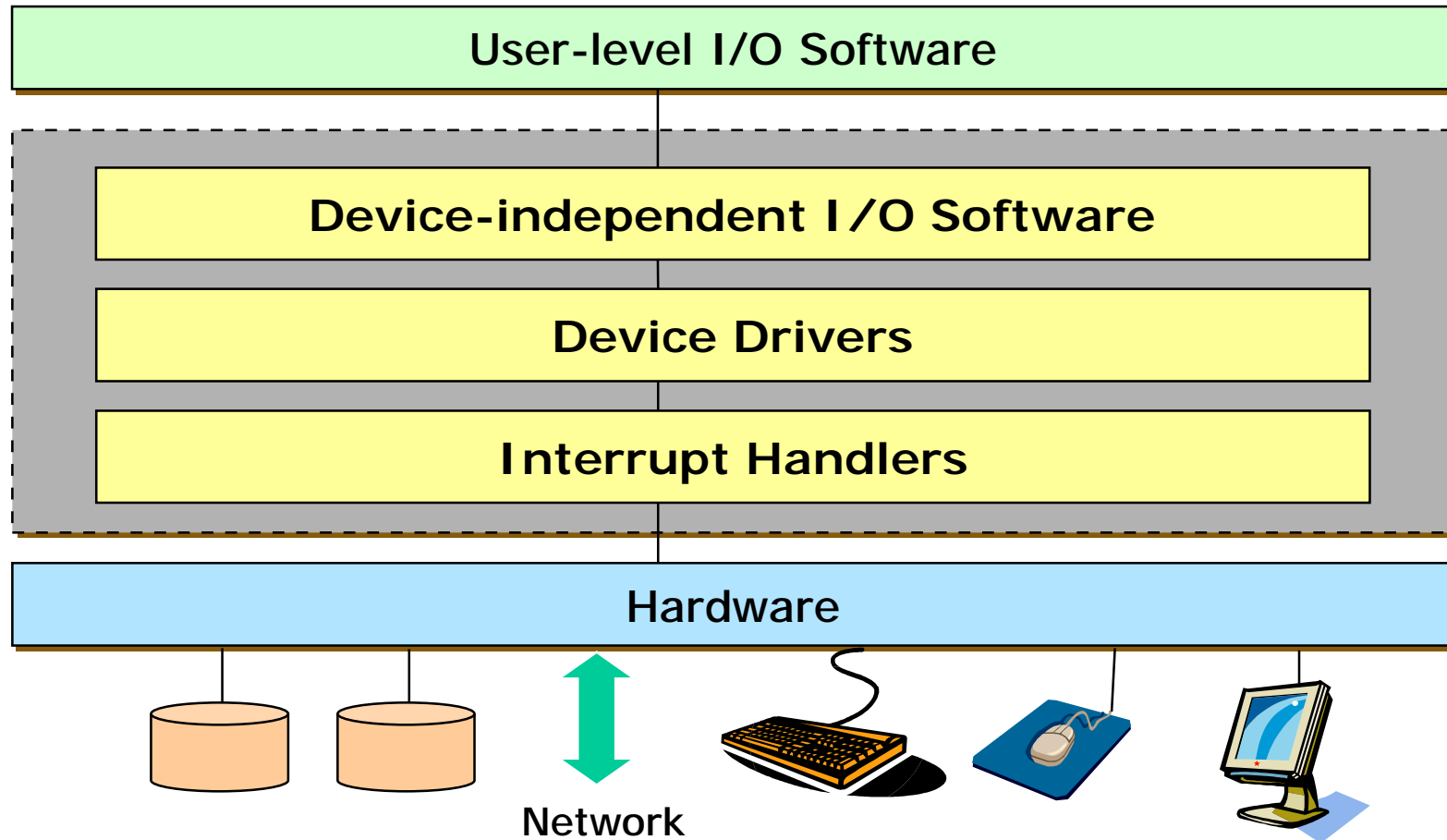- ü Buffering
    - § Data coming off a device cannot be stored in final destination
- ü Sharable vs. dedicated devices
    - § Disks vs. tape drives
    - § Unsharable devices introduce problems such as deadlocks

| User-level I/O Software |
|---|

| Device-independent I/O Software |
|---|
| Device Drivers |
| Interrupt Handlers |

| Hardware |
|---|

Network

**n** Handling interrupts

| Critical actions | : Acknowledge an interrupt to the PIC. |
|---|---|
|  | : Reprogram the PIC or the device controller. |
|  | : Update data structures accessed by both the device and the processor. |

**Reenable interrupts**

| Noncritical actions | : Update data structures that are accessed only by the processor. (e.g., reading the scan code from the keyboard) |
|---|---|

**Return from interrupts**

| Noncritical deferred actions | : Actions may be delayed. |
|---|---|
|  | : Copy buffer contents into the address space of some process (e.g., sending the keyboard line buffer to the terminal handler process). |
|  | **Bottom half (Linux)** |

**n** Device drivers

   **ü** Device-specific code to control each I/O device interacting with device-independent I/O software and interrupt handlers

   **ü** Requires to define a well-defined model and a standard interface of how they interact with the rest of the OS

   **ü** Implementing device drivers:

      **§** Statically linked with the kernel

      **§** Selectively loaded into the system during boot time

      **§** Dynamically loaded into the system during execution (especially for hot pluggable devices)

# Device-Independent I/O Software

**n** Uniform interfacing for device drivers

  **ü** In Unix, devices are modeled as special files

  § They are accessed through the use of system calls such as open(), read(), write(), close(), ioctl(), etc.

  § A file name is associated with each device

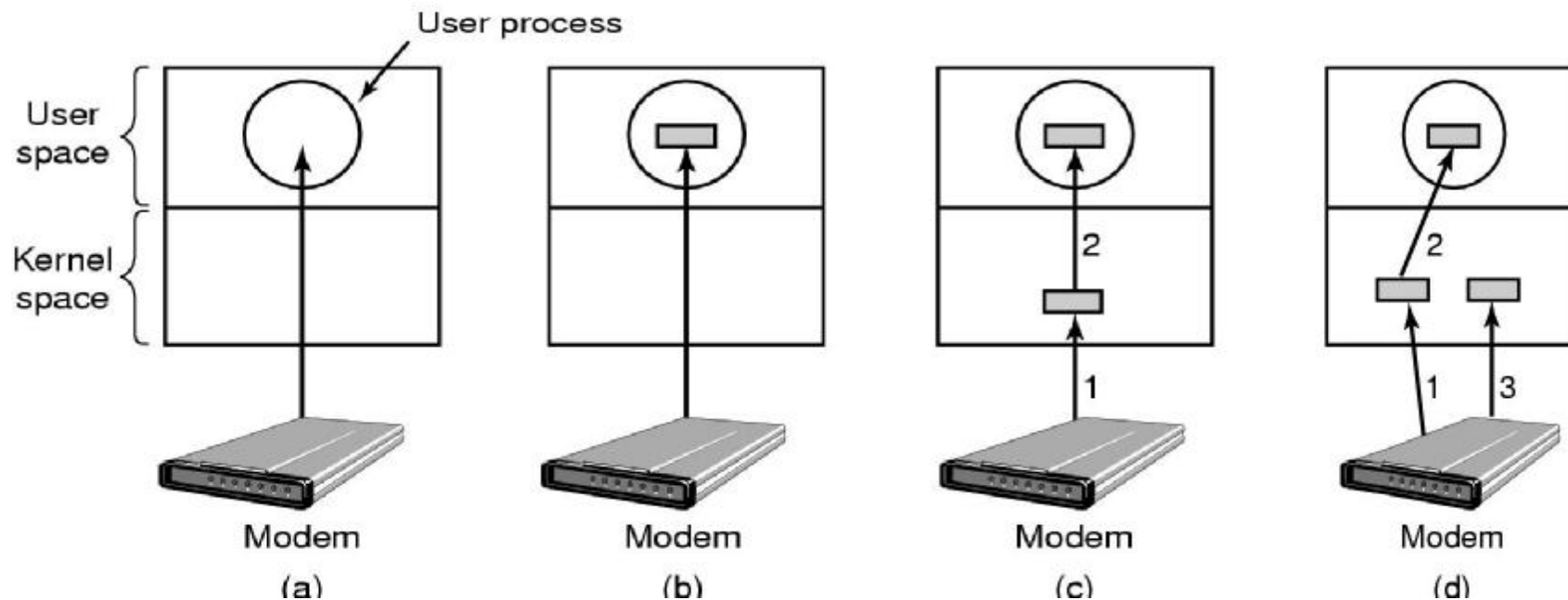  **ü** Major device number locates the appropriate driver

  § Minor device number (stored in i-node) is passed as a parameter to the driver in order to specify the unit to be read or written

  **ü** The usual protection rules for files also apply to I/O devices

**n** Buffering

   ü  (a) Unbuffered

   ü  (b) Buffered in user space

   ü  (c) Buffered in the kernel space

   ü  (d) Double buffering in the kernel

# Device-Independent I/O Software (Cont'd)

**n** Error reporting

    ü  Many errors are device-specific and must be handled by the appropriate driver, but the framework for error handling is device independent

    ü  Programming errors vs. actual I/O errors

    ü  Handling errors

        §  Returning the system call with an error code

        §  Retrying a certain number of times

        §  Ignoring the error

        §  Killing the calling process

        §  Terminating the system

# Device-Independent I/O Software (Cont'd)

**n  Allocating and releasing dedicated devices**

- ü  Some devices cannot be shared

- (1) Require processes to perform open()'s on the special files for devices directly

    - §  The process retries if open() fails

- (2) Have special mechanisms for requesting and releasing dedicated devices

    - §  An attempt to acquire a device that is not available blocks the caller

**n  Device-independent block size**

- ü  Treat several sectors as a single logical block

- ü  The higher layers only deal with abstract devices that all use the same block size

**n** Provided as a library

  ü Standard I/O library in C

    § fopen() vs. open()

**n** Spooling

  ü A way of dealing with dedicated I/O devices in a multiprogramming system

  ü Implemented by a daemon and a spooling directory

  ü Printers, network file transfers, USENET news, mails, etc.

| Layer | I/O functions |
|---|---|
| User processes | Make I/O call; format I/O; spooling |
| Device-independent software | Naming, protection, blocking, buffering, allocation |
| Device drivers | Set up device registers; check status |
| Interrupt handlers | Wake up driver when I/O completed |
| Hardware | Perform I/O operation |

I/O request

I/O reply