# Lecture 20
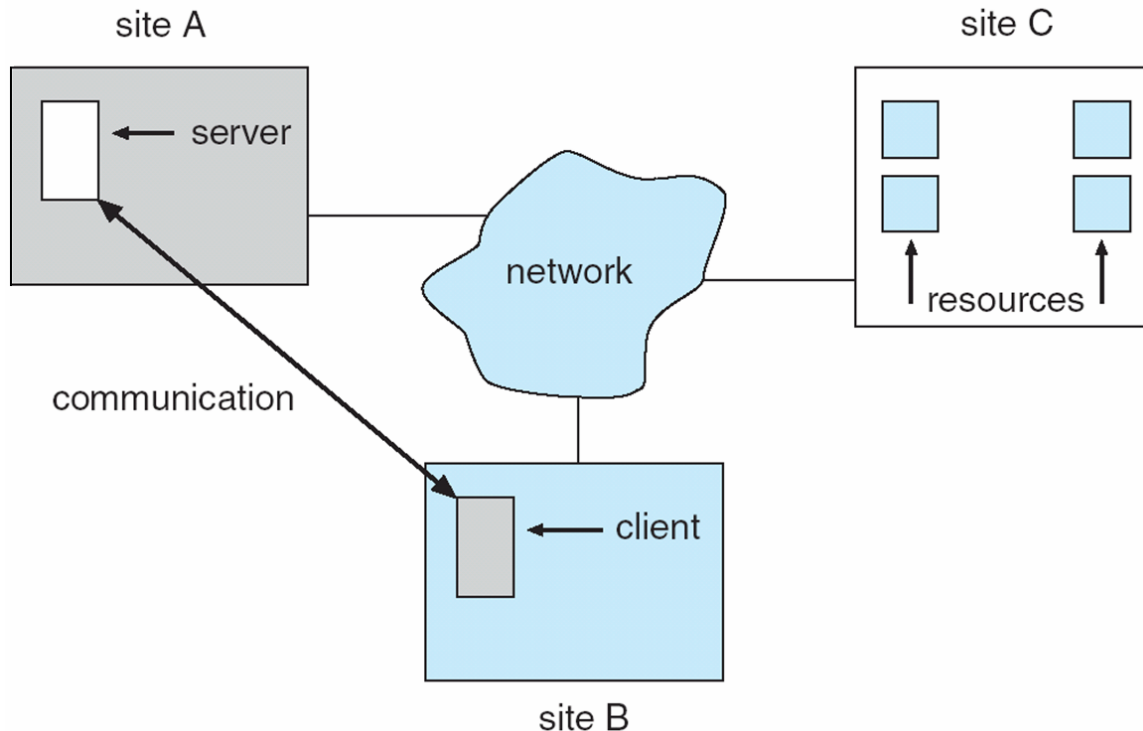
# Distributed System

- Distributed system is collection of loosely coupled processors interconnected by a communication network

- Processors variously called *nodes, computers, machines, hosts*
  - *Site* is location of the processor

# Motivation

Reasons for distributed systems

- Resource sharing

- Computation speedup – load sharing

- Reliability – detect and recover from site failure, function transfer, reintegrate failed site

- Communication – message passing

# Network-Operating Systems

- Users are aware of multiplicity of machines.  Access to resources of various machines is done explicitly by:

  - Remote logging into the appropriate remote machine (telnet, ssh)

  - Transferring data from remote machines to local machines, via the File Transfer Protocol (FTP) mechanism

# Distributed-Operating Systems

- Users not aware of multiplicity of machines. Access to remote resources similar to access to local resources

- Data Migration – transfer data by transferring entire file, or transferring only those portions of the file necessary for the immediate task

- Computation Migration – transfer the computation, rather than the data, across the system

# Distributed-Operating Systems (Cont.)

- Process Migration – execute an entire process, or parts of it, at different sites

  - **Load balancing** – distribute processes across network to even the workload

  - **Computation speedup** – subprocesses can run concurrently on different sites

  - **Hardware preference** – process execution may require specialized processor

  - **Software preference** – required software may be available at only a particular site

  - **Data access** – run process remotely, rather than transfer all data locally

# Example: Data Partioning

- Is a data migration scheme a good idea for computation speed-up for this code. If yes, give the data migration scheme.

```
for( i=1; i<=N; i++) {
    A[i] = A[i-1]*B[i];
}
```

# Design Issues

- **Fault tolerance** – the distributed system should continue to function in the face of failure

- **Scalability** – as demands increase, the system should easily accept the addition of new resources to accommodate the increased demand

- **Clusters** – a collection of semi-autonomous machines that acts as a single system

# Distributed File Systems

- Distributed file system (DFS) – a distributed implementation of the classical time-sharing model of a file system, where multiple users share files and storage resources

- A DFS manages set of dispersed storage devices

# DFS Structure

- **Service** – software entity running on one or more machines and providing a particular type of function to a priori unknown clients

- **Server** – service software running on a single machine

- **Client** – process that can invoke a service using a set of operations that forms its client interface

- A client interface for a file service is formed by a set of primitive file operations (create, delete, read, write)

- Client interface of a DFS transparent, i.e., not distinguish between local and remote files

# Naming Schemes — Three Main Approaches

- Files named by combination of their host name and local name; guarantees a unique system wide name

- Attach remote directories to local directories, giving the appearance of a coherent directory tree; only previously mounted remote directories can be accessed transparently

- Total integration of the component file systems
  - A single global name structure spans all the files in the system
  - If a server is unavailable, some arbitrary set of directories on different machines also becomes unavailable
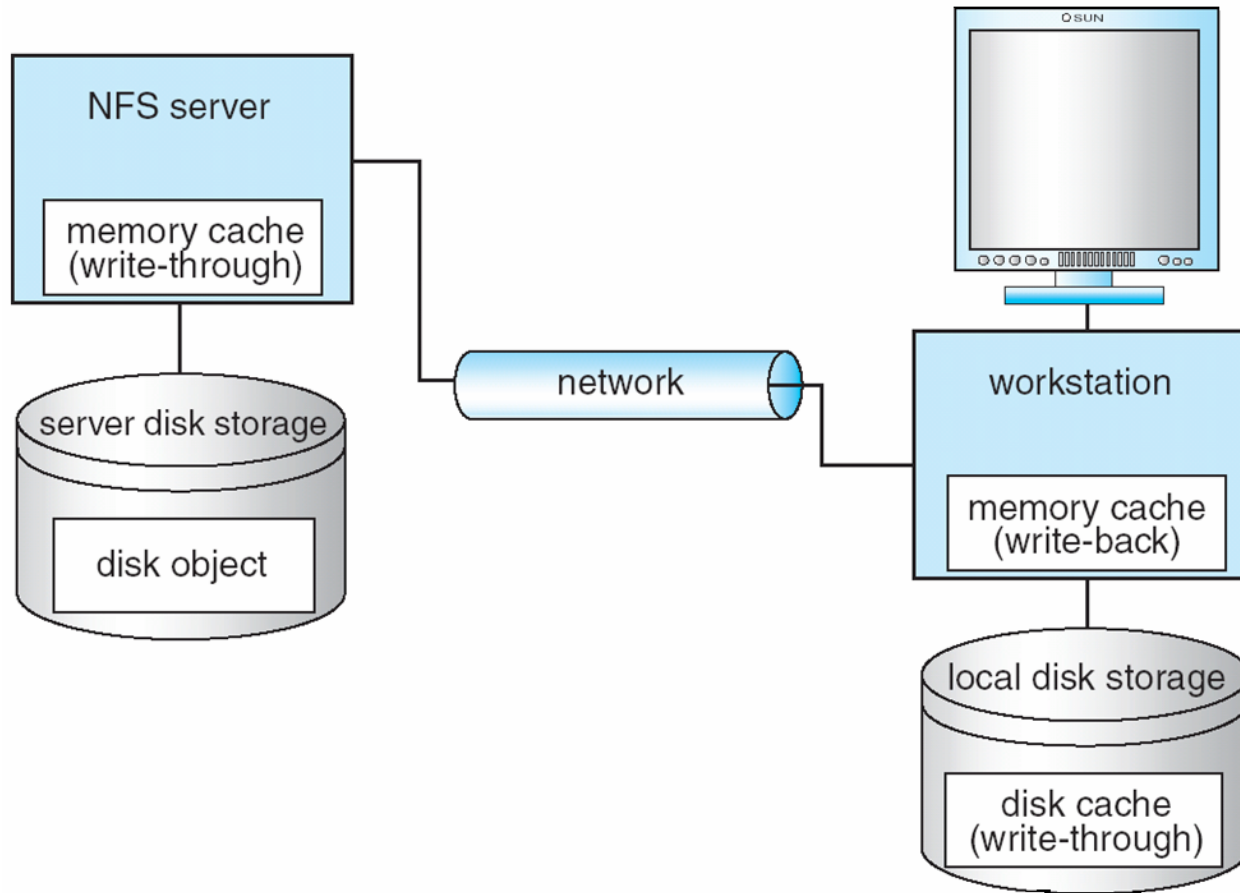
# Remote File Access

Reduce network traffic by retaining recently accessed disk blocks in a cache, so that repeated accesses to the same information can be handled locally

- If needed data not already cached, a copy of data is brought from the server to the user

- Accesses are performed on the cached copy

- Files identified with one master copy residing at the server machine, but copies of (parts of) the file are scattered in different caches

- Cache-consistency problem – keeping the cached copies consistent with the master file

  ‣ Could be called network virtual memory

# Cache Update Policy

- **Write-through** – write data through to disk as soon as they are placed on any cache
  - Reliable, but poor performance

- **Delayed-write** – modifications written to the cache and then written through to the server later
  - Write accesses complete quickly; some data may be overwritten before they are written back, and so need never be written at all
  - Poor reliability; unwritten data will be lost whenever a user machine crashes
  - Variation – scan cache at regular intervals and flush blocks that have been modified since the last scan
  - Variation – **write-on-close**, writes data back to the server when the file is closed
    - Best for files that are open for long periods and frequently modified

# Cachefs and its Use of Caching

# Remote Procedure Call

# Remote Procedure Call (RPC)

- Procedure call to another host

- Implementation:
  - Network protocol allow one procedure to call another procedure loaded on a different machine
  - Pass copy of parameters

# Network File System (NFS)

# The Sun Network File System (NFS)

- An implementation and a specification of a software system for accessing remote files across networks

- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP protocol and Ethernet

# NFS

■ A remote directory is mounted over a local file system directory

  ● The mounted directory looks like an integral  subtree of the local file system, replacing the subtree descending from the local directory

■ Specification of the remote directory for the mount operation is nontransparent; the host name of the remote directory has to be provided

  ● Files in the remote directory can then be accessed in a transparent manner
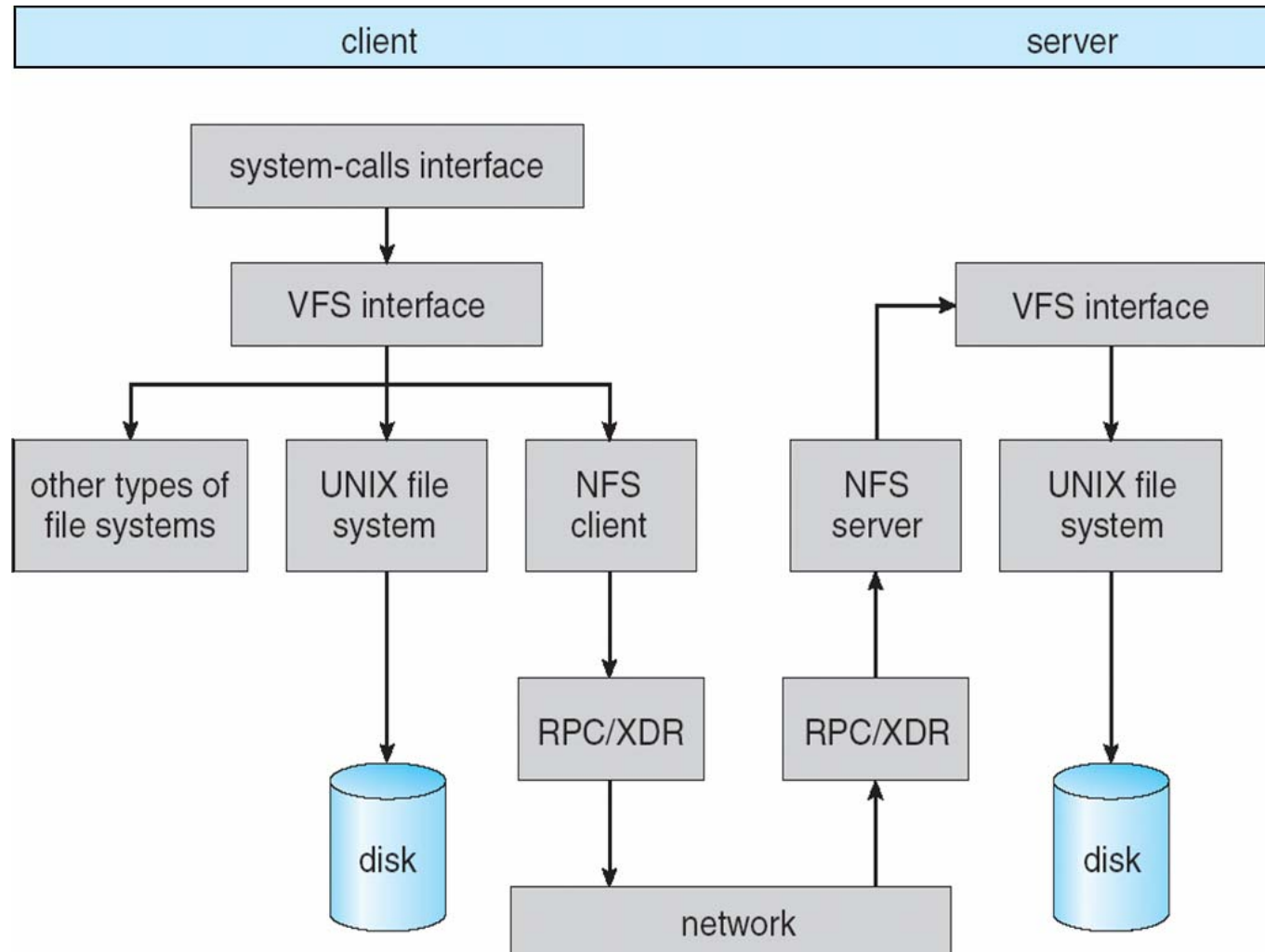
# NFS Protocol

- Provides a set of remote procedure calls for remote file operations.  The procedures support the following operations:
  - searching for a file within a directory
  - reading a set of directory entries
  - manipulating links and directories
  - accessing file attributes
  - reading and writing files
- NFS servers are **stateless**; each request has to provide a full set of arguments
  (NFS V4 is just coming available – very different, stateful)

# Three Major Layers of NFS Architecture

- UNIX file-system interface (based on the **open, read, write**, and **close** calls, and **file descriptors**)

- *Virtual File System* (VFS) layer – distinguishes local files from remote ones, and local files are further distinguished according to their file-system types
  - The VFS activates file-system-specific operations to handle local requests according to their file-system types
  - Calls the NFS protocol procedures for remote requests

- NFS service layer – bottom layer of the architecture
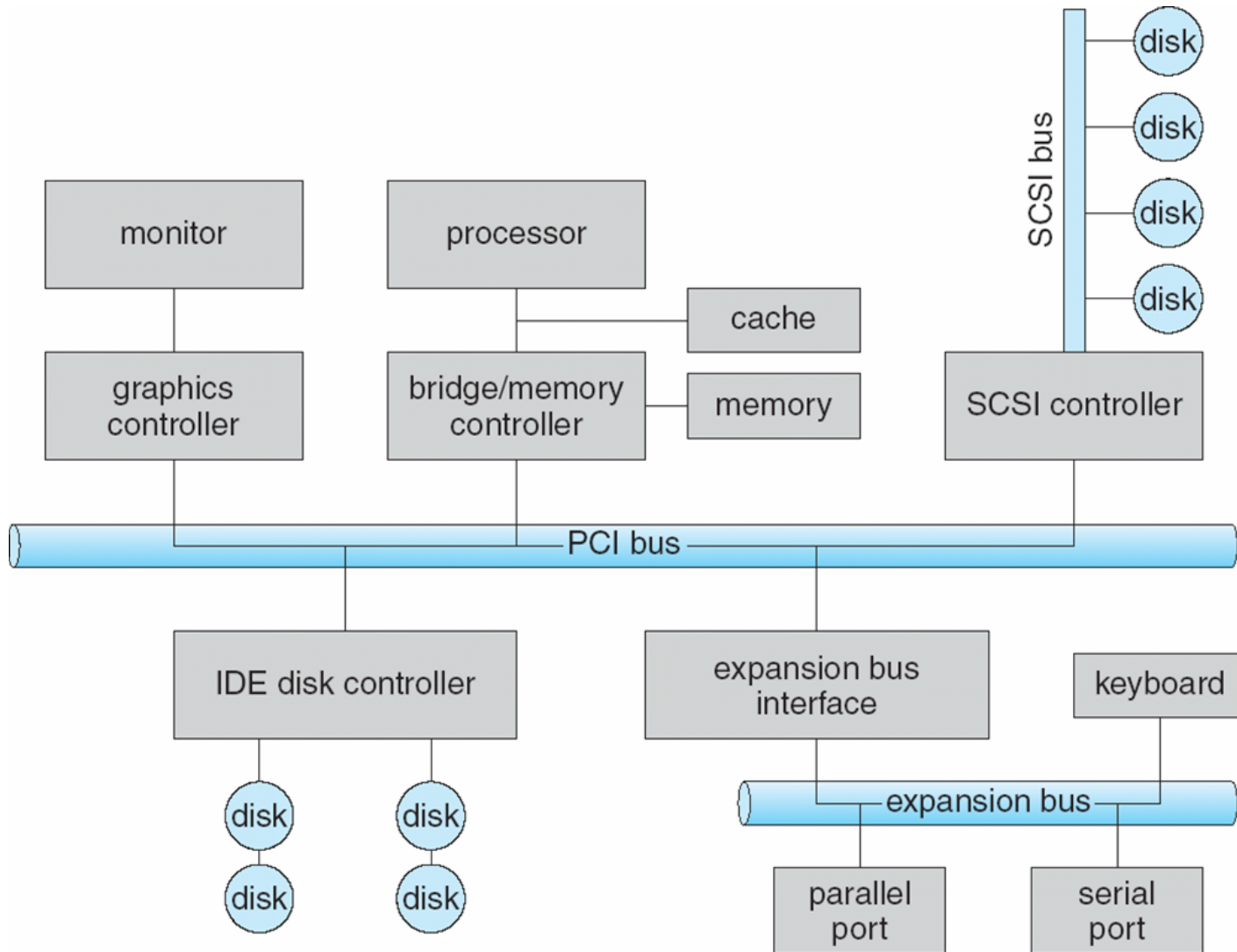  - Implements the NFS protocol

# Schematic View of NFS Architecture

# I/O Systems

# I/O Hardware

- Incredible variety of I/O devices

- Common concepts
  - Port
  - Bus
  - Controller

- I/O instructions control devices
- Devices have addresses, used by
  - Direct I/O instructions
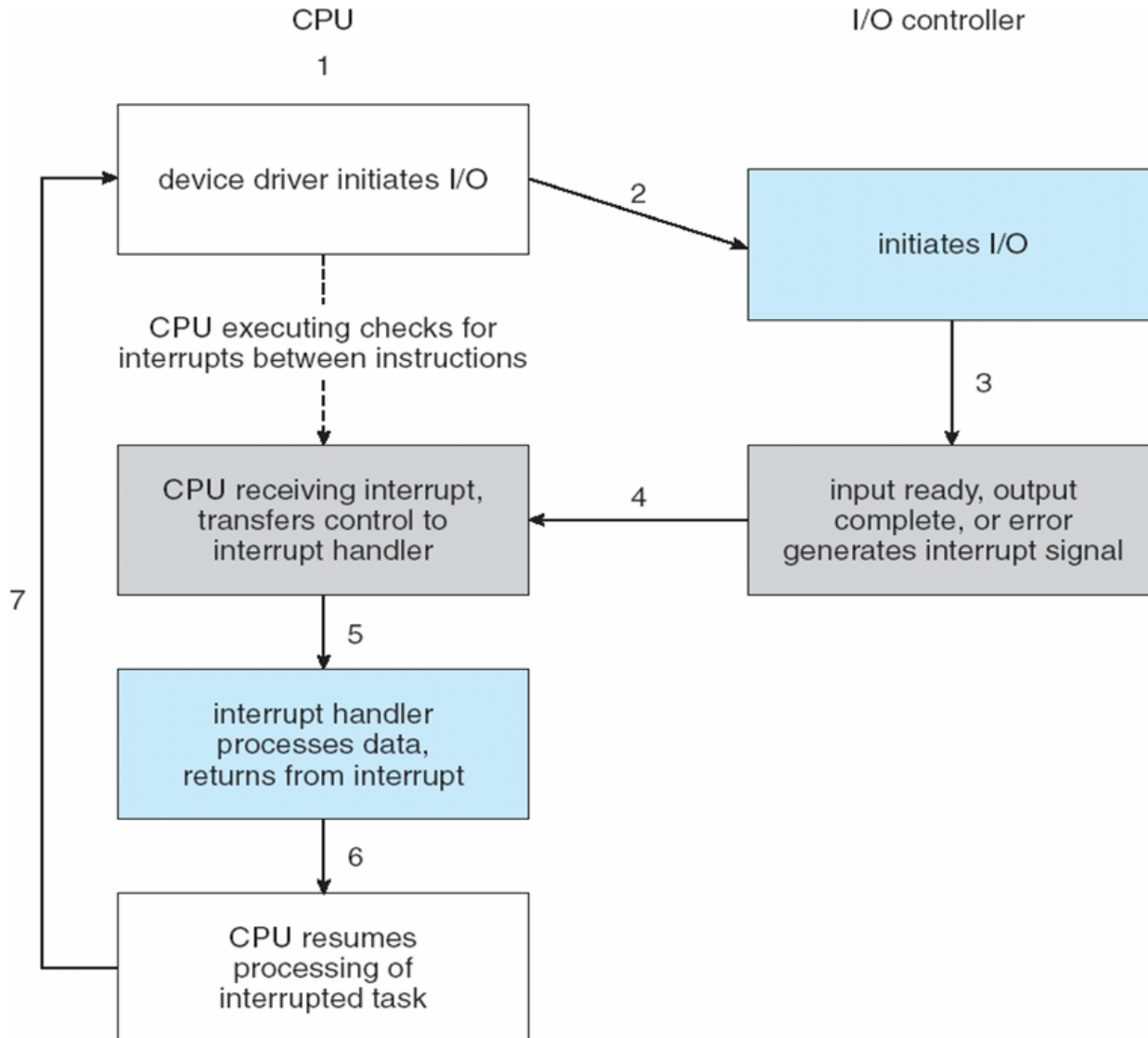  - Memory-mapped I/O

# A Typical PC Bus Structure

# Polling

- Determines state of device
    - command-ready
    - busy
    - Error
- Busy-wait cycle to wait for I/O from device
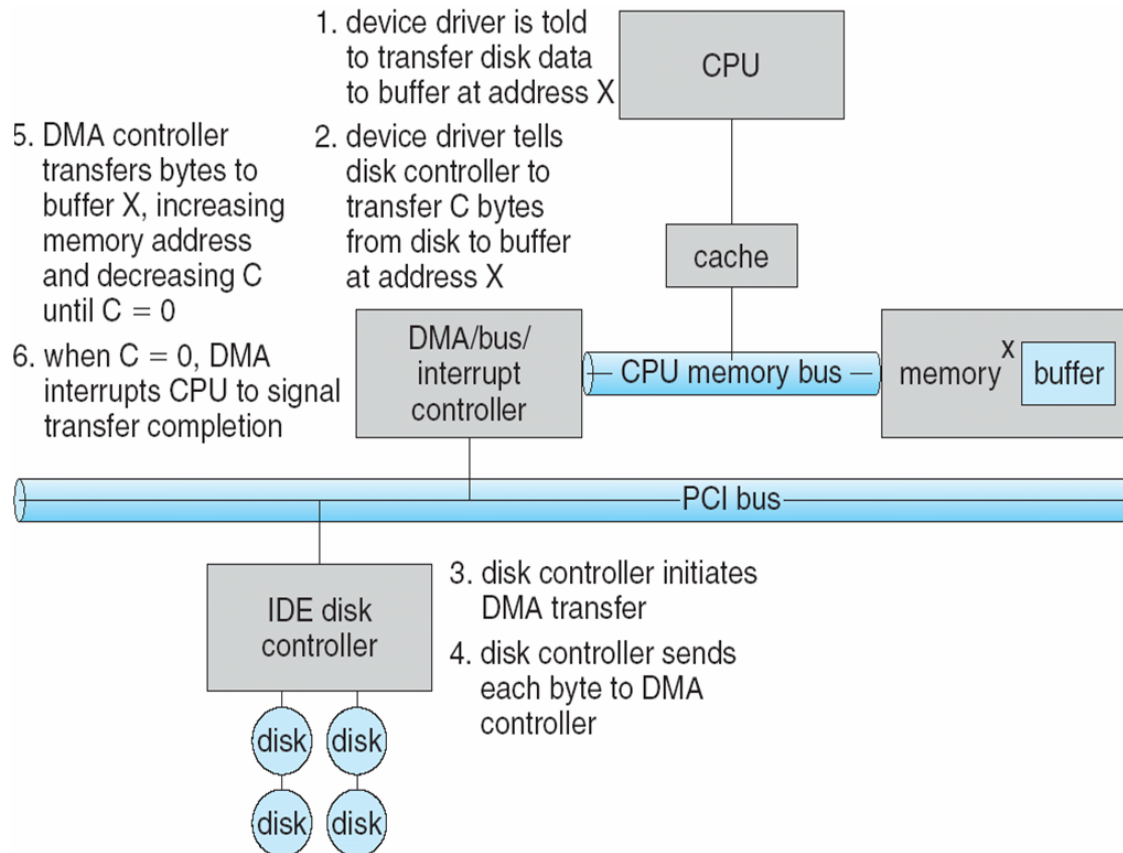
# Common Functions of Interrupts

- Interrupt transfers control to the interrupt service routine (ISR) generally, through the **interrupt vector**, addresses of all the service routines

- Interrupt handler preserves the state of the CPU by storing registers, the program counter, address of the interrupted instruction

- Incoming interrupts are *disabled* while another interrupt is being processed to prevent a *lost interrupt*

- A *trap* is a software-generated interrupt caused either by an error or a user request

- An operating system is **interrupt driven**
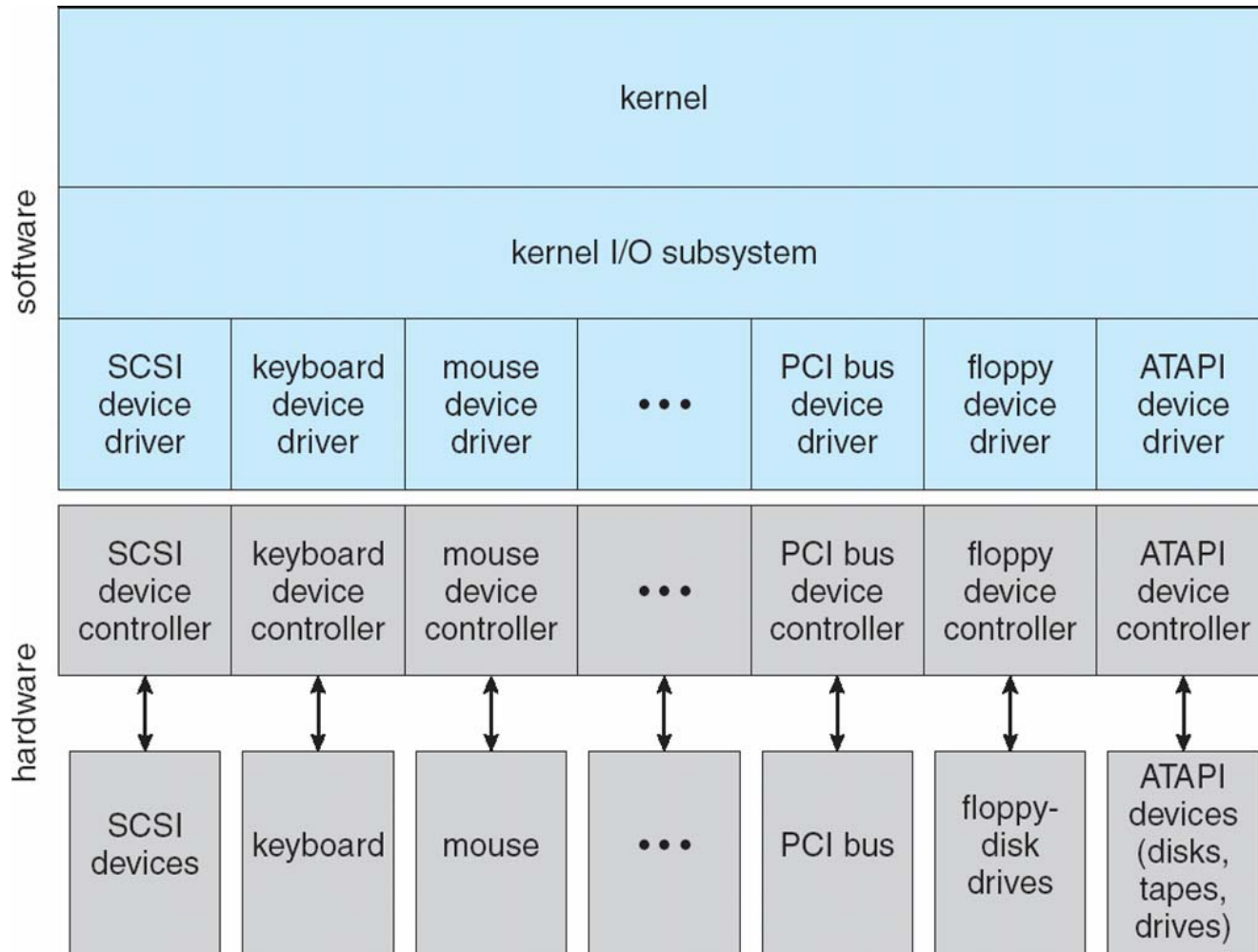
# Interrupt-Driven I/O Cycle

# Direct Memory Access

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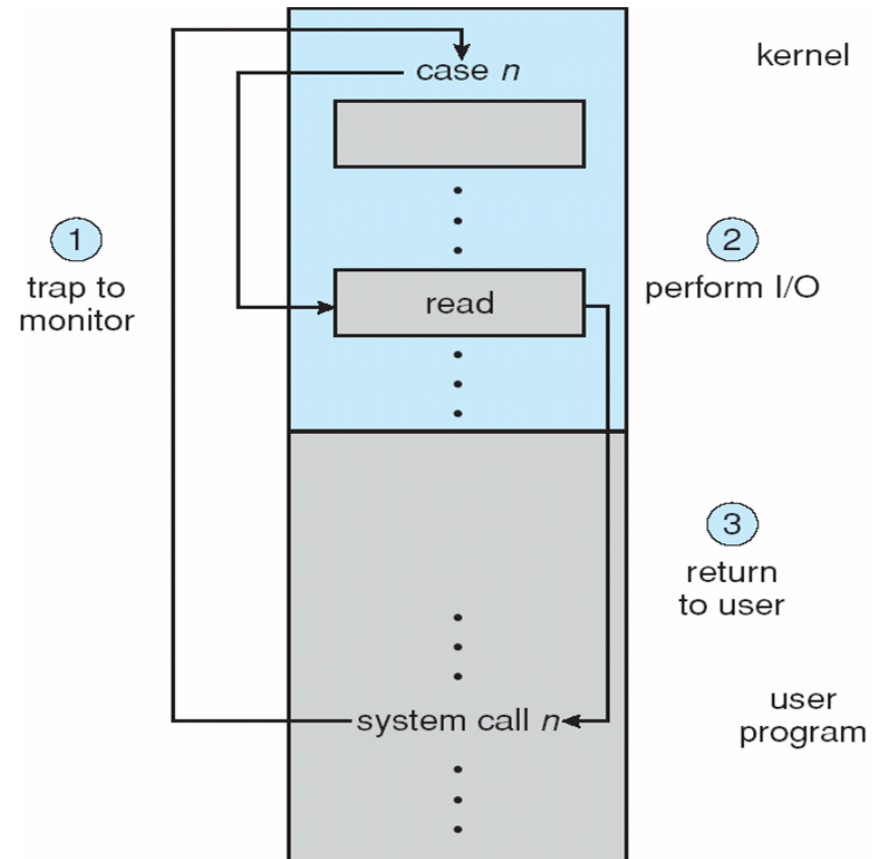- Bypasses CPU to transfer data directly between I/O device and memory

# A Kernel I/O Structure

# 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

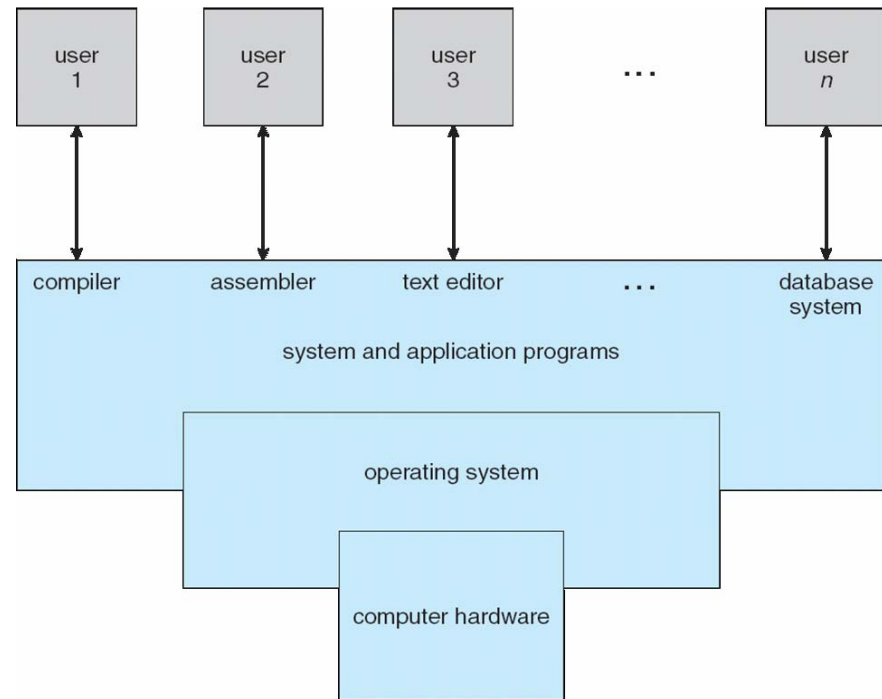    ▸ Memory-mapped and I/O port memory locations must be protected too

# Overview of Our Operating Systems Course Topics

# What is an Operating System?

- A program that acts as an intermediary between a user of a computer and the computer hardware

- Operating system goals:

    - Execute user programs and make solving user problems easier

    - Make the computer system convenient to use

    - Use the computer hardware in an efficient manner

# Computer System Structure

- Hardware – provides basic computing resources (CPU, memory, I/O devices)

- Operating system -- Controls and coordinates use of hardware among various applications and users

- Application programs – define the ways in which the system resources are used to solve the computing problems of the users (e.g. compilers, web browsers, database systems, video games)

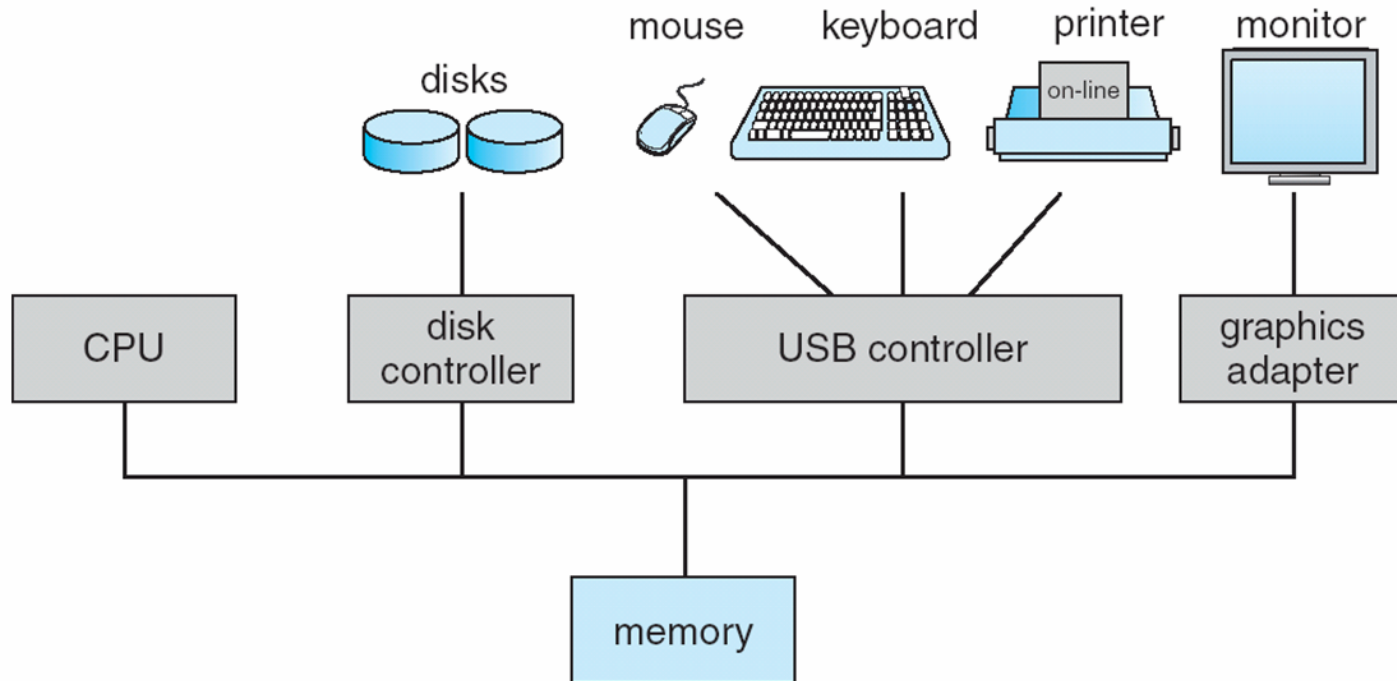- Users (e.g. People, machines, other computers)

# Computer Startup

- **bootstrap program** is loaded at power-up or reboot
  - Typically stored in ROM or EPROM, generally known as **firmware**
  - Initializes all aspects of system
  - Loads operating system kernel and starts execution

# Computer System Organization

- Computer-system operation
  - One or more CPUs, device controllers connect through common bus providing access to shared memory
  - Concurrent execution of CPUs and devices competing for memory cycles

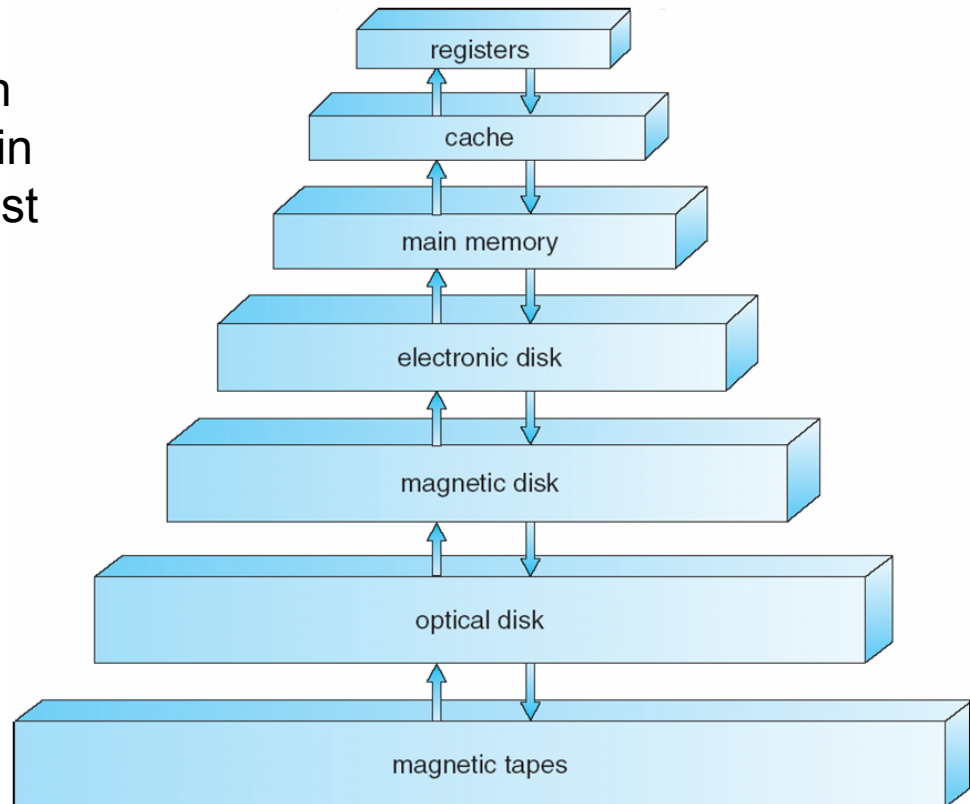# Computer-System Operation

- I/O devices and the CPU can execute concurrently

- Each device controller is in charge of a particular device type

- Each device controller has a local buffer

- CPU moves data from/to main memory to/from local buffers

- I/O is from the device to local buffer of controller

- Device controller informs CPU that it has finished its operation by causing an *interrupt*

# Storage Structure

- Main memory – only large storage media that the CPU can access directly

- Secondary storage – extension of main memory that provides large nonvolatile storage capacity

- Magnetic disks – rigid metal or glass platters covered with magnetic recording material

  - Disk surface is logically divided into **tracks**, which are subdivided into **sectors**

  - The **disk controller** determines the logical interaction between the device and the computer
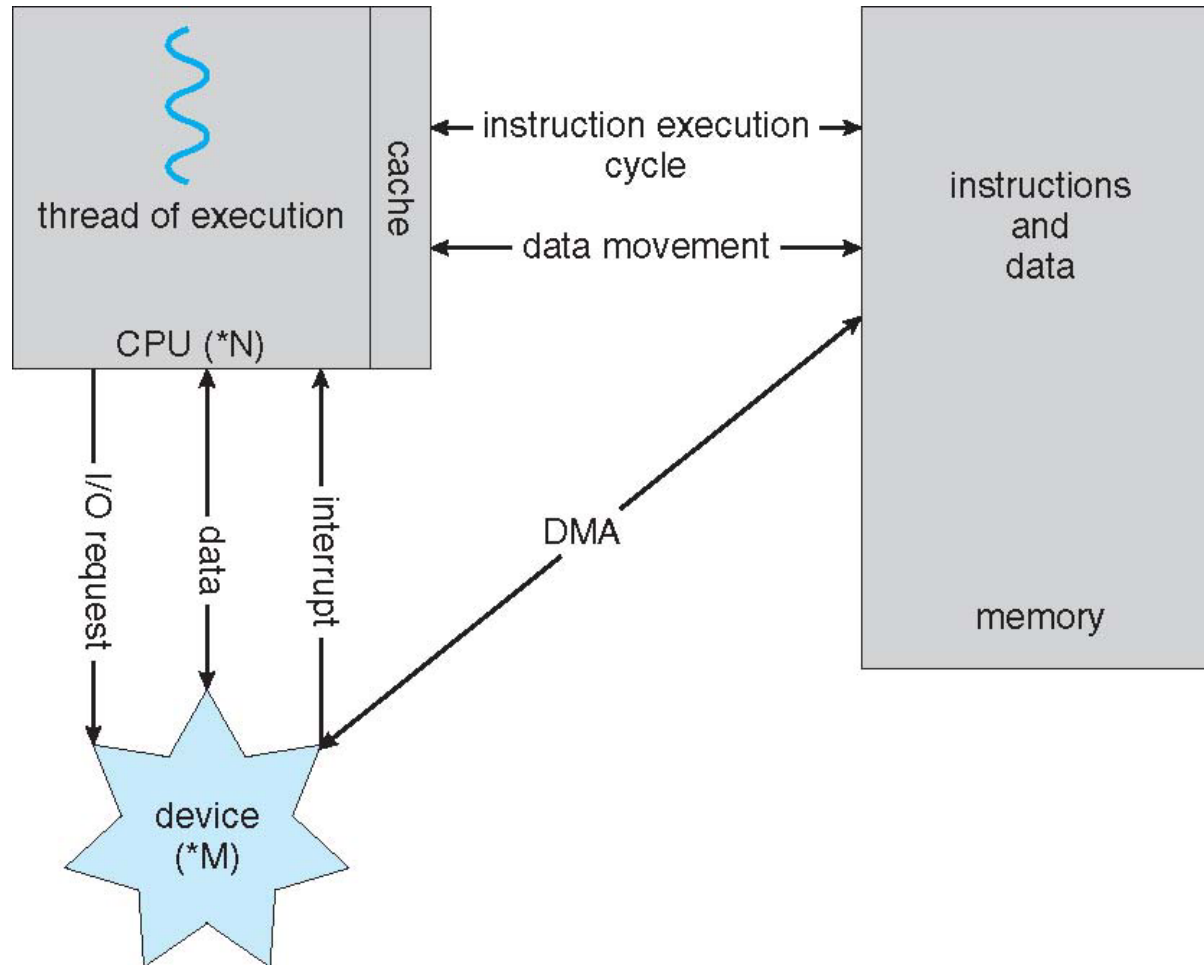
# Storage Hierarchy

- Storage systems organized in hierarchy
    - Speed
    - Cost
    - Volatility
- **Caching** – copying information into faster storage system; main memory can be viewed as a last *cache* for secondary storage

# Caching

- Important principle, performed at many levels in a computer (in hardware, operating system, software)

- Information in use copied from slower to faster storage temporarily

- Faster storage (cache) checked first to determine if information is there

  - If it is, information used directly from the cache (fast)

  - If not, data copied to cache and used there

- Cache smaller than storage being cached

  - Cache management important design problem

  - Cache size and replacement policy

# How a Modern Computer Works

# A Dual-Core Design

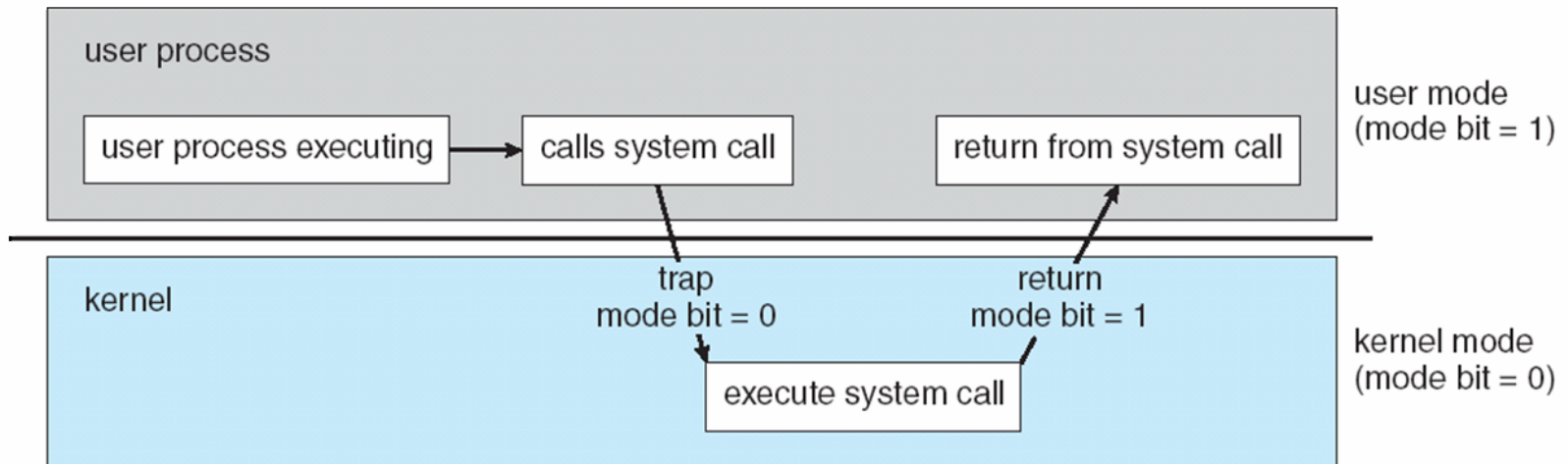# Operating System Structure

- **Multiprogramming** needed for efficiency
  - Single user cannot keep CPU and I/O devices busy at all times
  - Multiprogramming organizes jobs (code and data) so CPU always has one to execute
  - A subset of total jobs in system is kept in memory
  - One job selected and run via **job scheduling**
  - When it has to wait (for I/O for example), OS switches to another job
- **Timesharing (multitasking)** is logical extension in which CPU switches jobs so frequently that users can interact with each job while it is running, creating **interactive** computing
  - **Response time** should be < 1 second
  - Each user has at least one program executing in memory ⇨ **process**
  - If several jobs ready to run at the same time ⇨ **CPU scheduling**
  - If processes don't fit in memory, **swapping** moves them in and out to run
  - **Virtual memory** allows execution of processes not completely in memory

# Operating-System Operations

- Interrupt driven by hardware

- Software error or request creates **exception** or **trap**
  - Division by zero, request for operating system service

- Other process problems include infinite loop, processes modifying each other or the operating system

- **Dual-mode** operation allows OS to protect itself and other system components
  - **User mode** and **kernel mode**
  - **Mode bit** provided by hardware
    - Provides ability to distinguish when system is running user code or kernel code
    - Some instructions designated as **privileged**, only executable in kernel mode
    - System call changes mode to kernel, return from call resets it to user

# Transition from User to Kernel Mode

■ Timer to prevent infinite loop / process hogging resources

- Set interrupt after specific period
- Operating system decrements counter
- When counter zero generate an interrupt
- Set up before scheduling process to regain control or terminate program that exceeds allotted time

# Process Management

- A process is a program in execution. It is a unit of work within the system. Program is a *passive entity*, process is an *active entity*.

- Process needs resources to accomplish its task
  - CPU, memory, I/O, files
  - Initialization data

- Process termination requires reclaim of any reusable resources

- Single-threaded process has one **program counter** specifying location of next instruction to execute
  - Process executes instructions sequentially, one at a time, until completion
  - Multi-threaded process has one program counter per thread

- Typically system has many processes, some user, some operating system running concurrently on one or more CPUs
  - Concurrency by multiplexing the CPUs among the processes / threads

# Process Management Activities

The operating system is responsible for the following activities in connection with process management:

- Creating and deleting both user and system processes

- Suspending and resuming processes

- Providing mechanisms for process synchronization

- Providing mechanisms for process communication

- Providing mechanisms for deadlock handling

# Memory Management

- All data in memory before and after processing

- All instructions in memory in order to execute

- Memory management determines what is in memory when
  - Optimizing CPU utilization and computer response to users

- Memory management activities
  - Keeping track of which parts of memory are currently being used and by whom
  - Deciding which processes (or parts thereof) and data to move into and out of memory
  - Allocating and deallocating memory space as needed

# Storage Management

- OS provides uniform, logical view of information storage
  - Abstracts physical properties to logical storage unit  - **file**

- File-System management
  - Files usually organized into directories
  - Access control on most systems to determine who can access what
  - OS activities include
    - Creating and deleting files and directories
    - Primitives to manipulate files and dirs
    - Mapping files onto secondary storage
    - Backup files onto stable (non-volatile) storage media

# Mass-Storage Management

- Usually disks used to store data that does not fit in main memory or data that must be kept for a "long" period of time

- Proper management is of central importance

- Entire speed of computer operation hinges on disk subsystem and its algorithms

- OS activities
  - Free-space management
  - Storage allocation
  - Disk scheduling

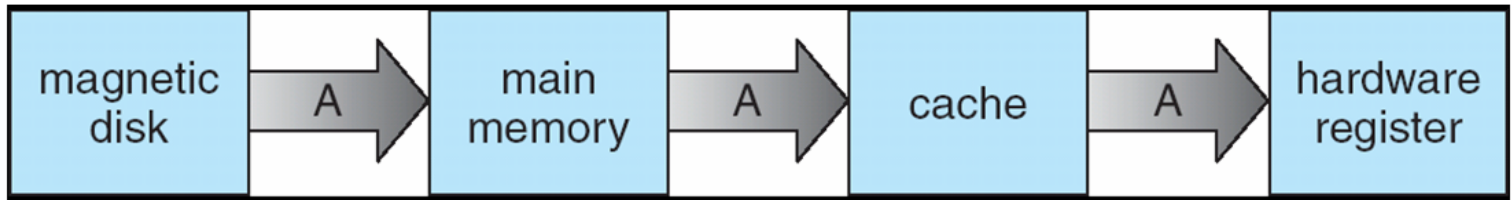- Tertiary storage includes optical storage, magnetic tape

# Performance of Various Levels of Storage

■ Movement between levels of storage hierarchy can be explicit or implicit

| Level | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Name | registers | cache | main memory | disk storage |
| Typical size | < 1 KB | > 16 MB | > 16 GB | > 100 GB |
| Implementation technology | custom memory with multiple ports, CMOS | on-chip or off-chip CMOS SRAM | CMOS DRAM | magnetic disk |
| Access time (ns) | 0.25 – 0.5 | 0.5 – 25 | 80 – 250 | 5,000.000 |
| Bandwidth (MB/sec) | 20,000 – 100,000 | 5000 – 10,000 | 1000 – 5000 | 20 – 150 |
| Managed by | compiler | hardware | operating system | operating system |
| Backed by | cache | main memory | disk | CD or tape |

# Migration of Integer A from Disk to Register

■ Multitasking environments must be careful to use most recent value, no matter where it is stored in the storage hierarchy

| magnetic disk | →A→ | main memory | →A→ | cache | →A→ | hardware register |

■ Multiprocessor environment must provide cache coherency in hardware such that all CPUs have the most recent value in their cache

■ Distributed environment situation even more complex
  ● Several copies of a datum can exist

# I/O Subsystem

- One purpose of OS is to hide peculiarities of hardware devices from the user

- I/O subsystem responsible for

  - Memory management of I/O including buffering (storing data temporarily while it is being transferred), caching (storing parts of data in faster storage for performance), spooling (the overlapping of output of one job with input of other jobs)

  - General device-driver interface

  - Drivers for specific hardware devices

# Protection and Security

- **Protection** – any mechanism for controlling access of processes or users to resources defined by the OS

- **Security** – defense of the system against internal and external attacks
  - Huge range, including denial-of-service, worms, viruses, identity theft, theft of service

- Systems generally first distinguish among users, to determine who can do what
  - **Privilege escalation** allows user to change to effective ID with more rights

# Open-Source Operating Systems

- Operating systems made available in source-code format rather than just binary closed-source

- Counter to the copy protection and Digital Rights Management (DRM) movement

- Started by Free Software Foundation (FSF), which has "copyleft" GNU Public License (GPL)

- Examples include GNU/Linux, BSD UNIX (including core of Mac OS X), and Sun Solaris

**See you at the final**