

Lecture 1

Introduction & Overview

Operating Systems



Objectives

- To provide functionalities of the major operating systems components
- To provide principles of modern operating systems. In particular, the course will cover details of concurrent processes, multi-threads, CPU scheduling, memory management, file system, storage subsystem, and input/output management



Contents

■ Process Scheduling

- processes and threads, context switching, synchronization, scheduling, and deadlock.

■ Memory Management

- linking, dynamic memory allocation, dynamic address translation, virtual memory, and demand paging.

■ File Systems

- storage devices, disk management and scheduling

Grading Scheme

■ Mid-Term	20%
■ Assignments	10%
■ Quizzes	5%
■ Semester Mini Project	5%
■ Final	60%
■ <i>Total</i>	<i>100 points</i>



Course Material

■ Text Book

- Abraham Silberschatz and Peter Baer Galvin:
Operating System Concepts, Addison-Wesley
(8th ed.)

■ Reference Book

- Andrew S. Tanenbaum and Albert S.
Woodhull, Operating Systems: Design and
Implementation, Prentice Hall (3rd ed.)



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

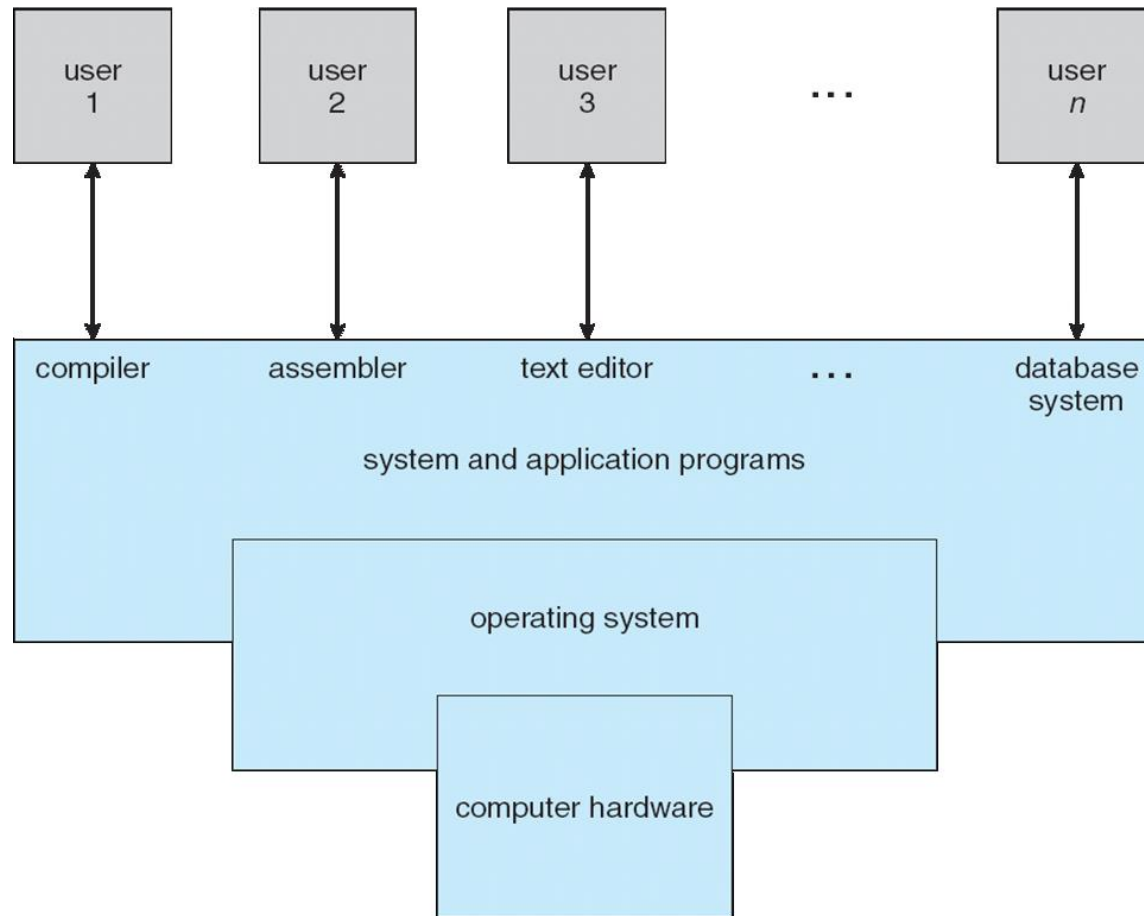


Computer System Structure

■ Computer system can be divided into four components:

- 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
 - Word processors, compilers, web browsers, database systems, video games
- Users
 - People, machines, other computers

Four Components of a Computer System



What Operating Systems Do?

- Depends on the point of view
- Users want convenience, **ease of use**
 - Don't care about **resource utilization**
- But shared computer such as **mainframe** or **minicomputer** must keep all users happy
- Users of dedicate systems such as **workstations** have dedicated resources but frequently use shared resources from **servers**
- Handheld computers are resource poor, optimized for usability and battery life
- Some computers have little or no user interface, such as embedded computers in devices and automobiles



Operating System Definition

- OS is a **resource allocator**

- ☐ Manages all resources
- ☐ Decides between conflicting requests for efficient and fair resource use

- OS is a **control program**

- ☐ Controls execution of programs to prevent errors and improper use of the computer

Operating System Definition (Cont.)

- No universally accepted definition
- “Everything a vendor ships when you order an operating system” is good approximation
 - But varies wildly
- “The one program running at all times on the computer” is the **kernel**.
Everything else is either a system program (ships with the operating system) or an application program.

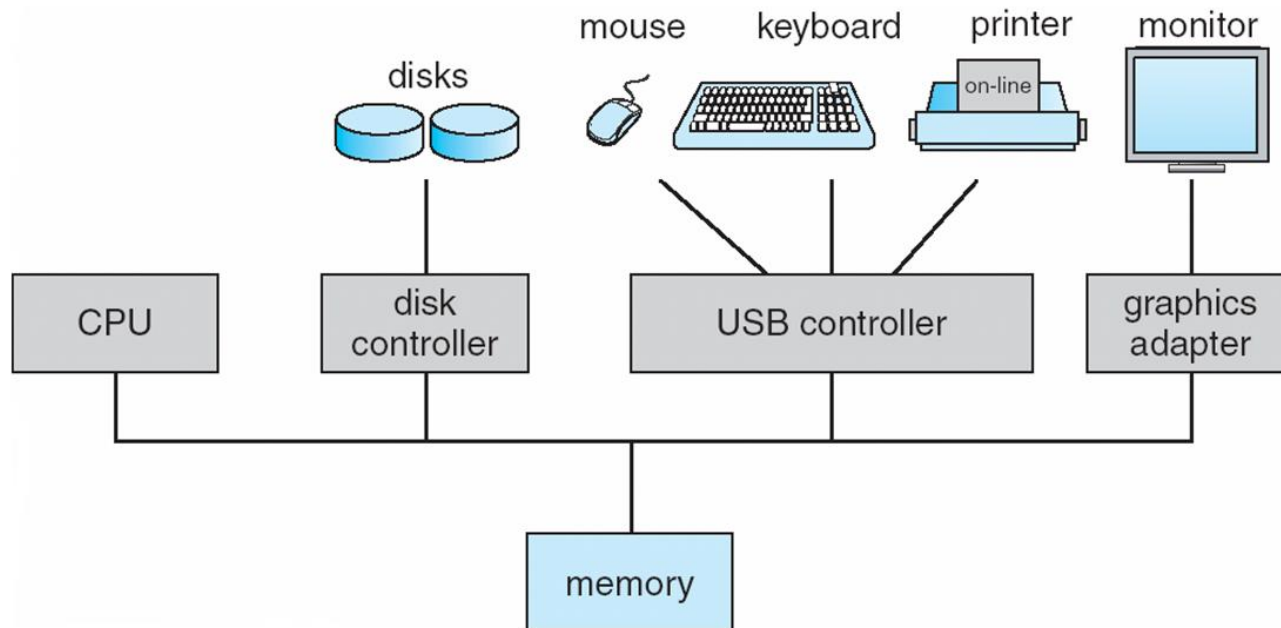
Computer Startup

- **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

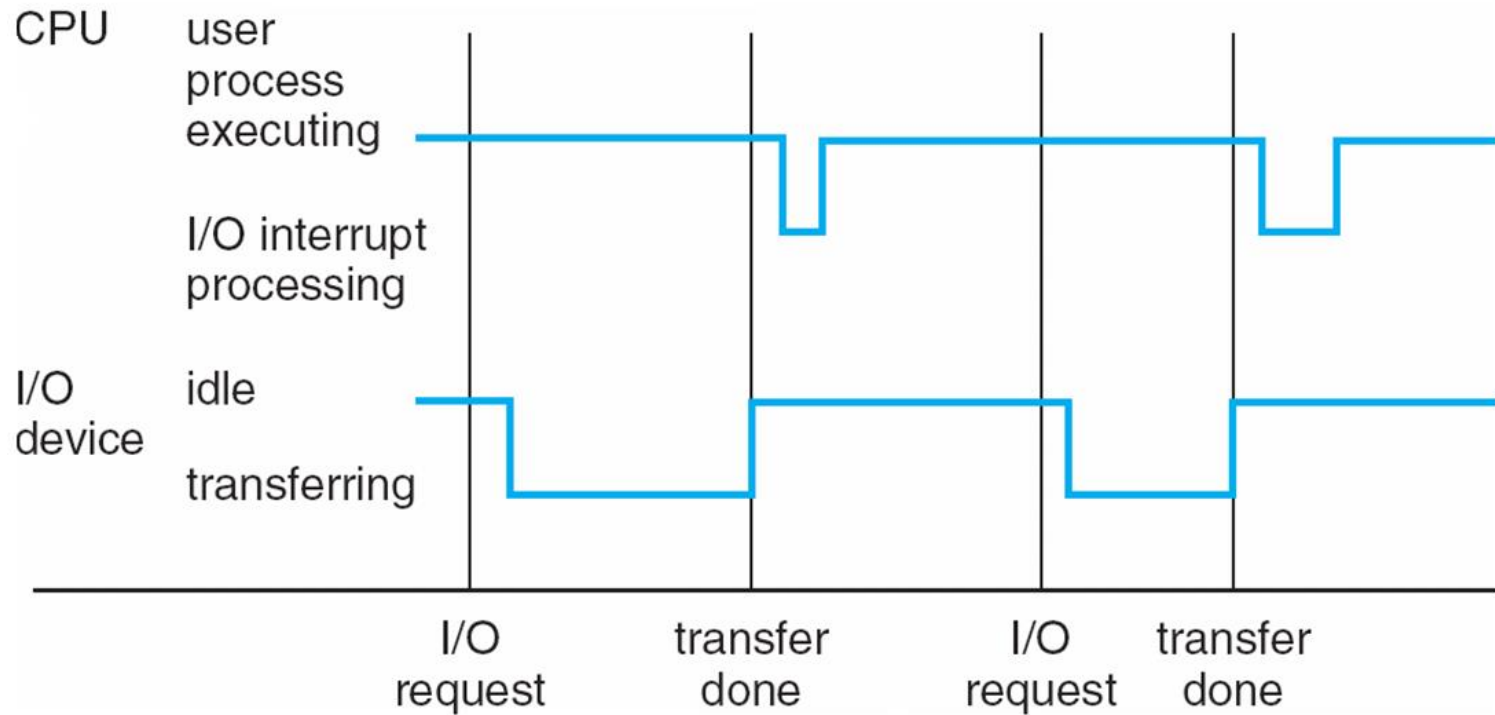




Computer-System Operation

- 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**

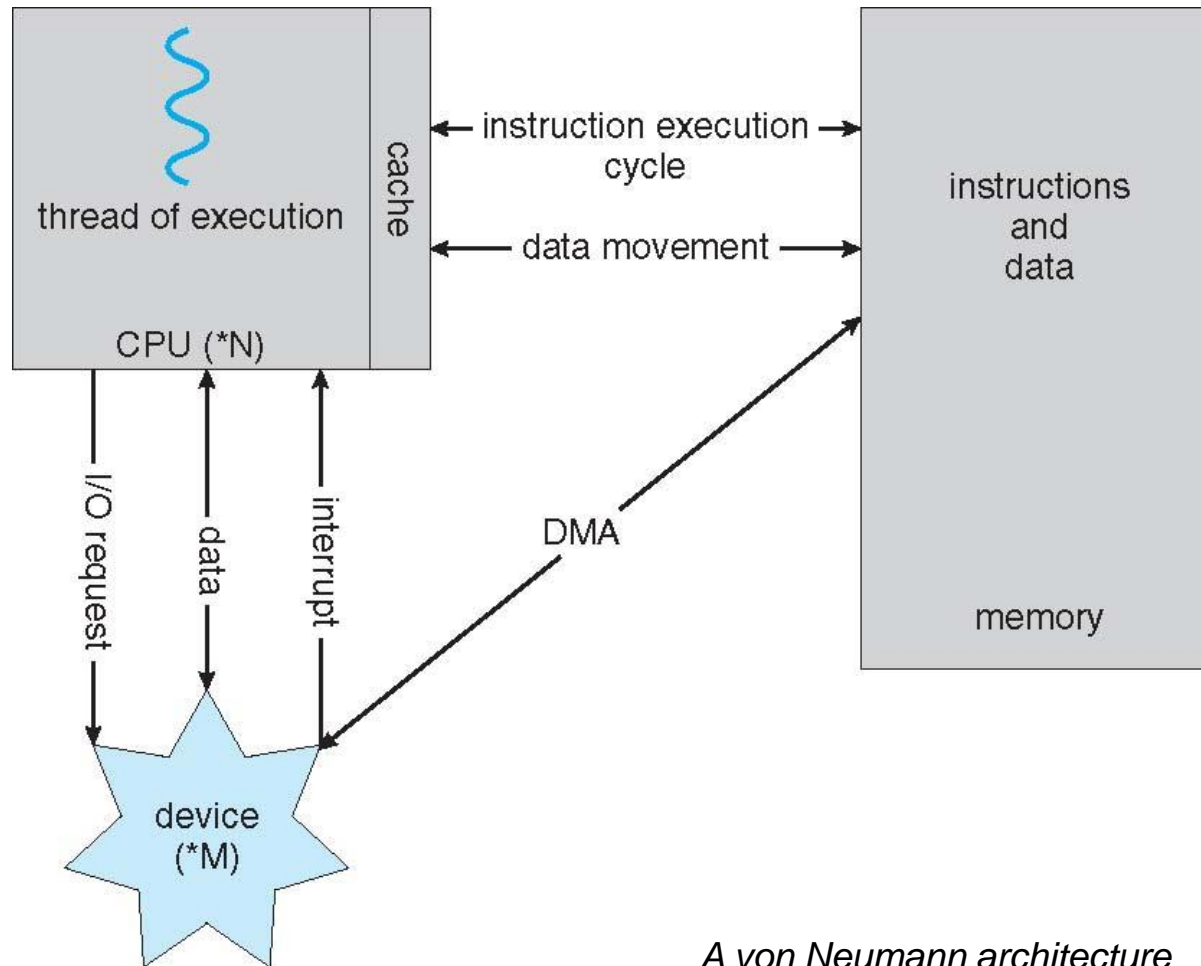
Interrupt Timeline



Computer-System Architecture

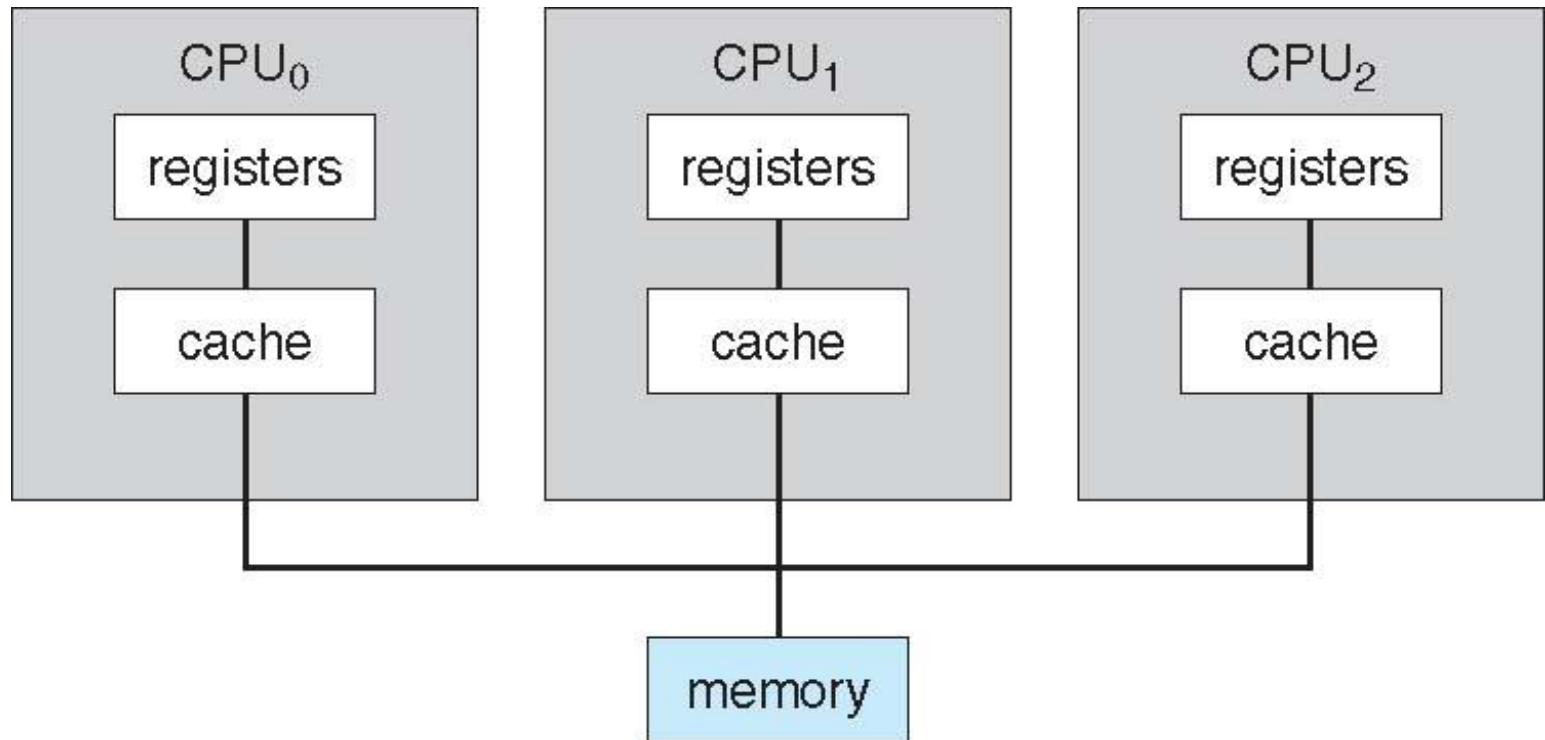
- Most systems use a single general-purpose processor (PDAs through mainframes)
 - Most systems have special-purpose processors as well
- **Multiprocessors** systems growing in use and importance
 - Also known as **parallel systems**, **tightly-coupled systems**
 - Advantages include:
 1. **Increased throughput**
 2. **Economy of scale**
 3. **Increased reliability – graceful degradation** or **fault tolerance**
 - Two types:
 1. **Asymmetric Multiprocessing**
 2. **Symmetric Multiprocessing**

How a Modern Computer Works

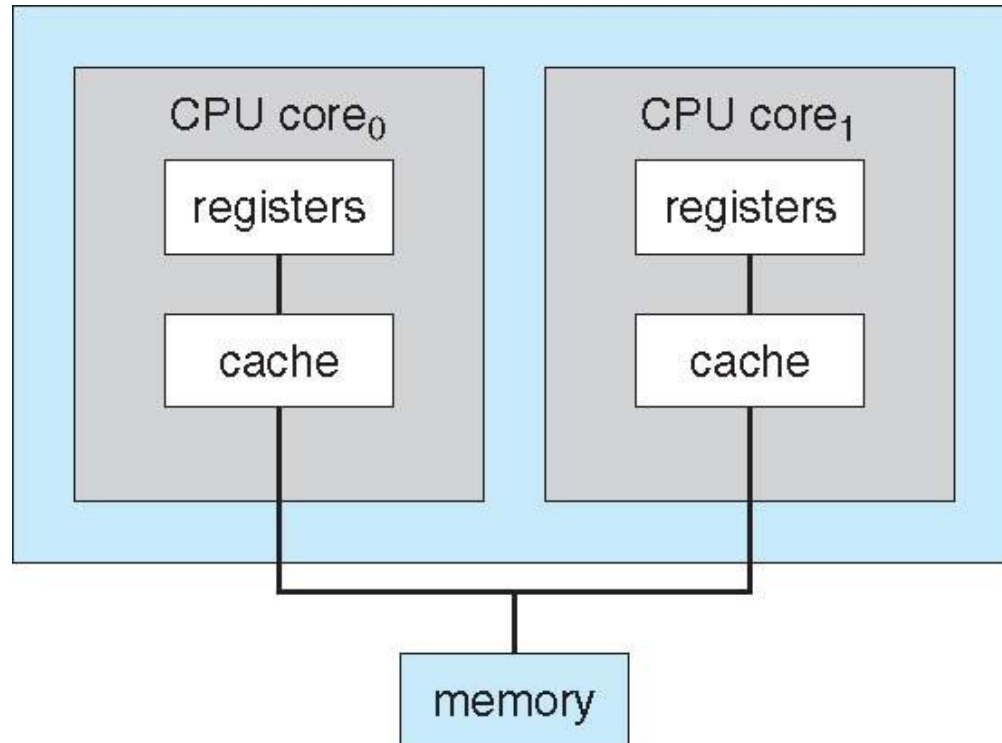


A von Neumann architecture

Symmetric Multiprocessing Architecture



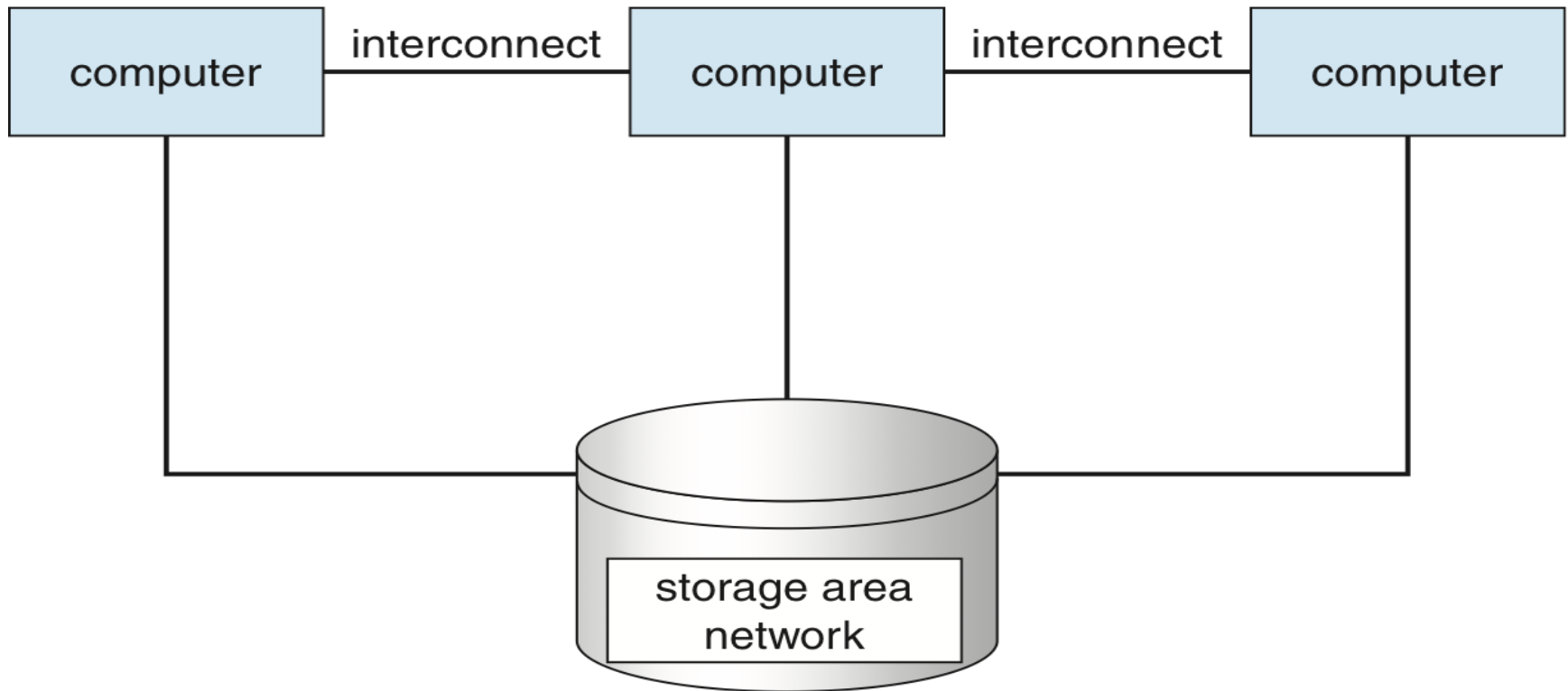
A Dual-Core Design



Clustered Systems

- Like multiprocessor systems, but multiple systems working together
 - Usually sharing storage via a **storage-area network (SAN)**
 - Provides a **high-availability** service which survives failures
 - **Asymmetric clustering** has one machine in hot-standby mode
 - **Symmetric clustering** has multiple nodes running applications, monitoring each other
 - Some clusters are for **high-performance computing (HPC)**
 - Applications must be written to use **parallelization**

Clustered Systems



Types of Operating System

Simple Batch System

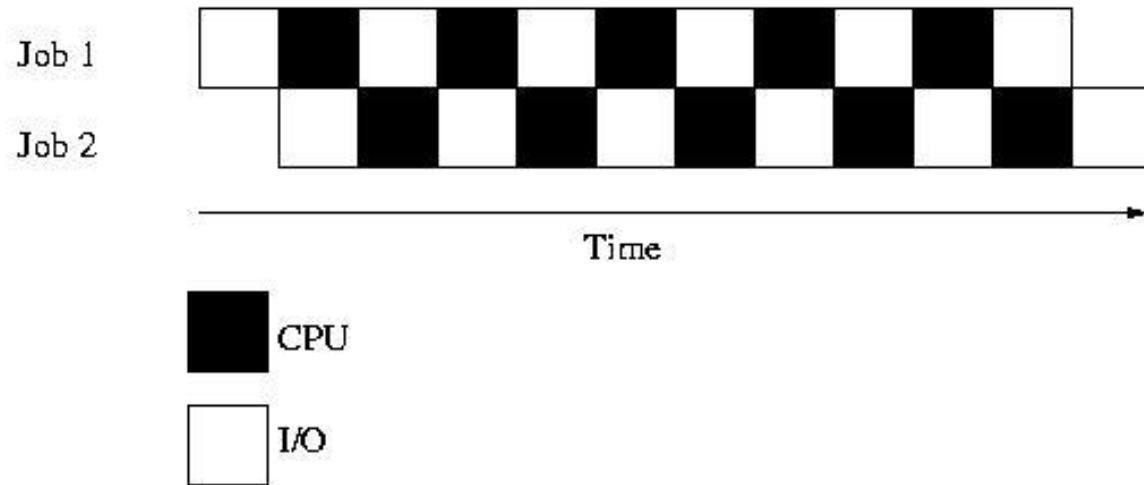
- The user submits a job (written on a card or tape) to a computer operator.
- The computer operator place a batch of several jobs on a input device
- A special program, the monitor, manages the execution of each program in the batch
- Resident monitor is in the main memory and available for execution

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

Multiprogramming

- One job can use the CPU while the other is waiting for I/O



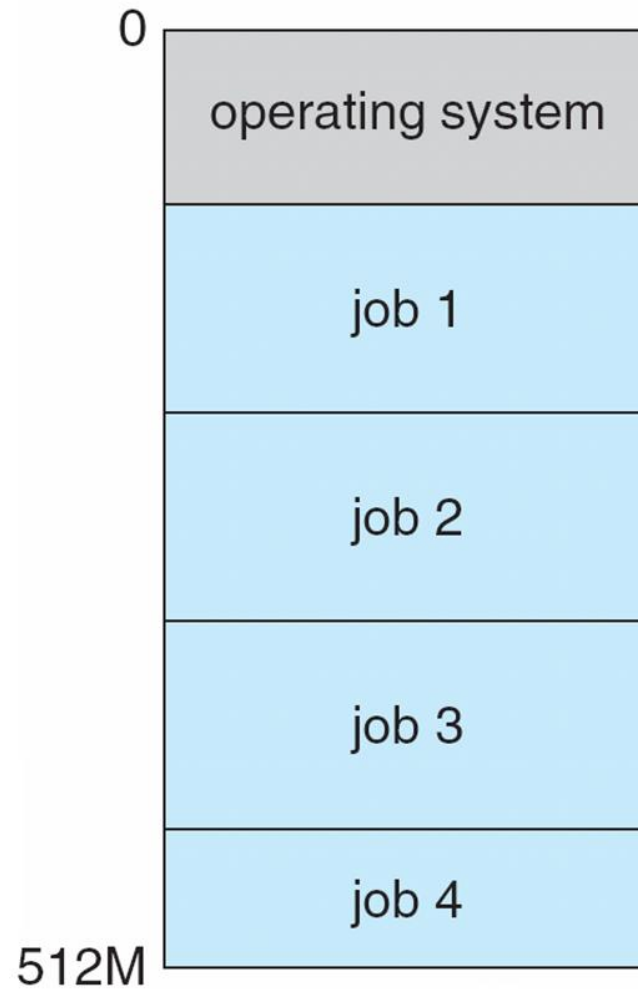
- Small jobs not delayed by large jobs
- Overhead?
- Context switching

Time Sharing System (TSS)

- Processor's time is shared among multiple users
- Use cheap terminals to let multiple users interact with the system at the same time.
- OS does timesharing to give illusion of each user has own computer
- User can interact during the execution time
- Plus Multiple jobs can be run
- User can
 - Interact and React
 - Control the path of the program
 - Perform interactive debugging

- **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 \Rightarrow **process**
 - If several jobs ready to run at the same time \Rightarrow **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

Memory Layout for Multiprogrammed System





Personal Computing

- Computers are cheap, so give everyone a computer.
- Initially, OS became a subroutine library again (MSDos)
- Since then, adding back in memory protection, multiprogramming, etc.
- Because when humans are expensive, don't waste their time by letting programs crash each other

Distributed Computing

- Collection of separate, possibly heterogeneous, systems networked together
 - Network is a communications path
 - Local Area Network (**LAN**)
 - Wide Area Network (**WAN**)
 - Metropolitan Area Network (**MAN**)
- Network Operating System provides features between systems across network
 - Communication scheme allows systems to exchange messages
 - Illusion of a single system

Special-Purpose Systems

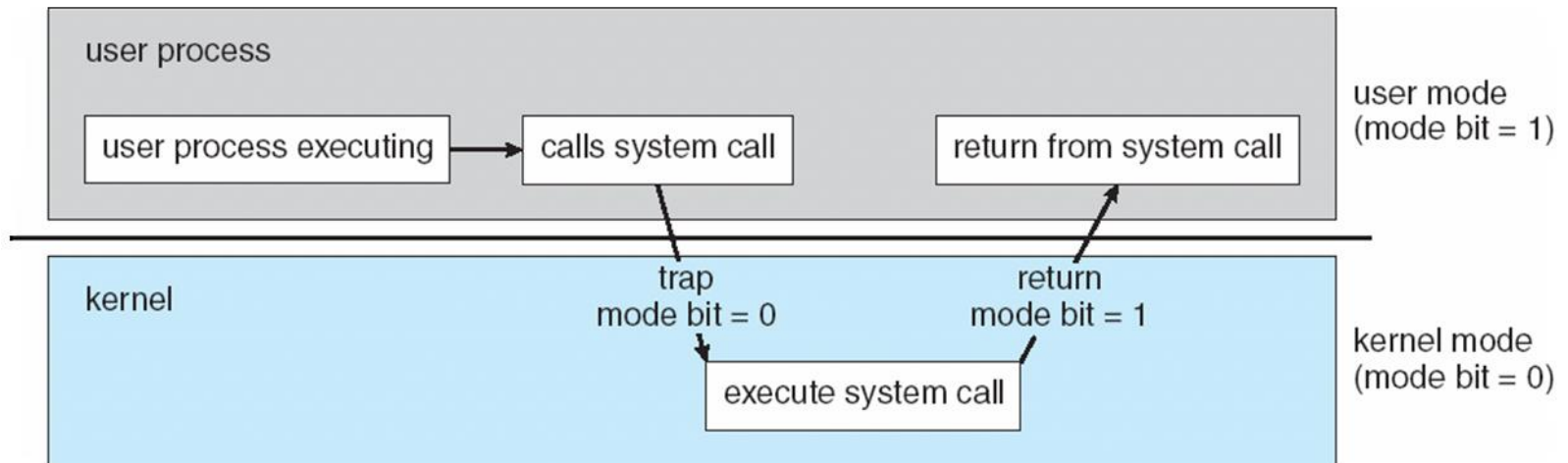
- Real-time embedded systems most prevalent form of computers
 - Vary considerable, special purpose, limited purpose OS, **real-time OS**
- Multimedia systems
 - Streams of data must be delivered according to time restrictions
- Handheld systems
 - PDAs, smart phones, limited CPU, memory, power
 - Reduced feature set OS, limited I/O

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**
 - Each medium is controlled by device (i.e., disk drive, tape drive)
 - Varying properties include access speed, capacity, data-transfer rate, access method (sequential or random)
- 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