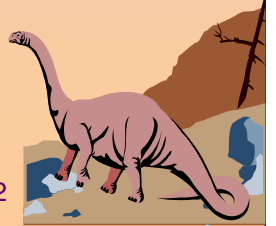# Chapter 1:   Introduction

- ■ What is an Operating System?
- ■ Mainframe Systems
- ■ Desktop Systems
- ■ Multiprocessor Systems
- ■ Distributed Systems
- ■ Clustered System
- ■ Real -Time Systems
- ■ Handheld Systems
- ■ Computing Environments

# 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.

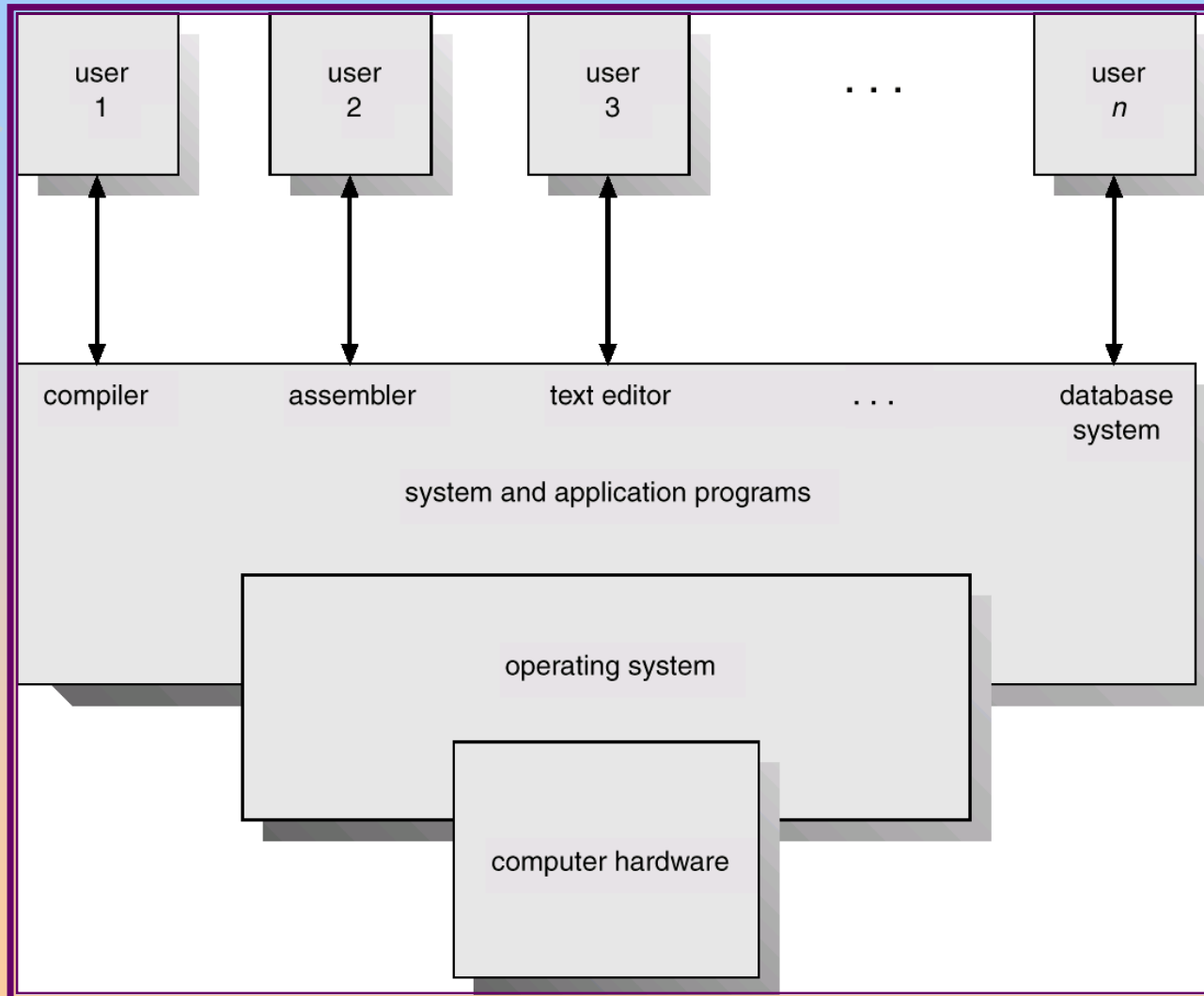- Use the computer hardware in an efficient manner.

# Computer System Components

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

2. Operating system – controls and coordinates the use of the hardware among the various application programs for the various users.

3. Applications programs – define the ways in which the system resources are used to solve the computing problems of the users (compilers, database systems, video games, business programs).

4. Users (people, machines, other computers).

# Abstract View of System Components



| user 1 | user 2 | user 3 | . . . | user n |

compiler     assembler     text editor     . . .     database system

system and application programs

operating system
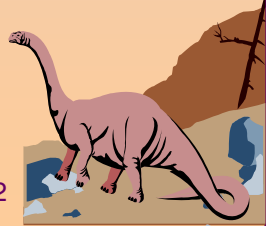
computer hardware

# Operating System Definitions

- Resource allocator – manages and allocates resources.
- Control program – controls the execution of user programs and operations of I/O devices .
- Kernel – the one program running at all times (all else being application programs).
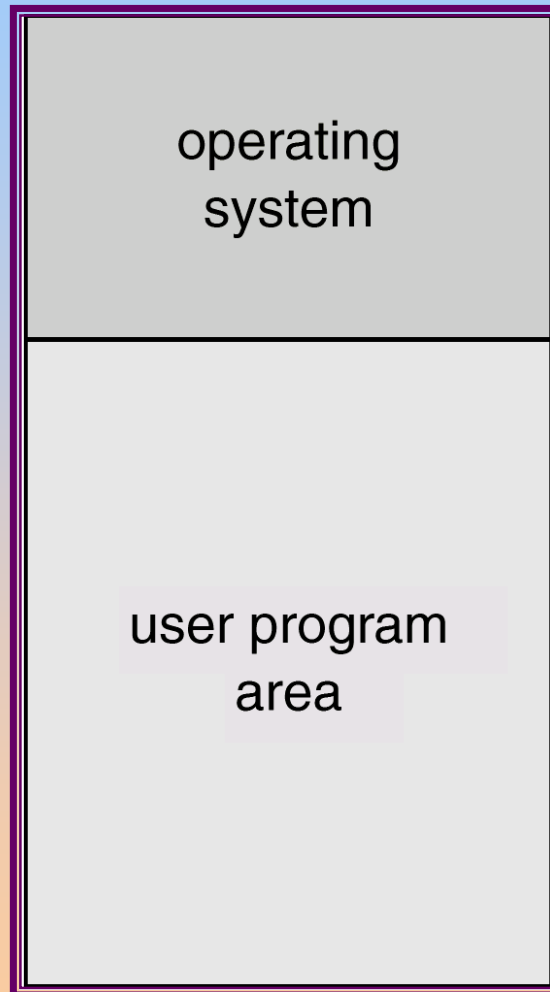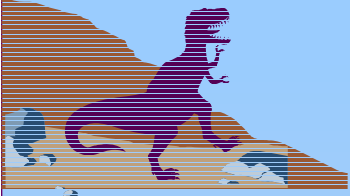
# Mainframe Systems

- Reduce setup time by batching similar jobs

- Automatic job sequencing – automatically transfers control from one job to another.  First rudimentary operating system.

- Resident monitor
  - initial control in monitor
  - control transfers to job
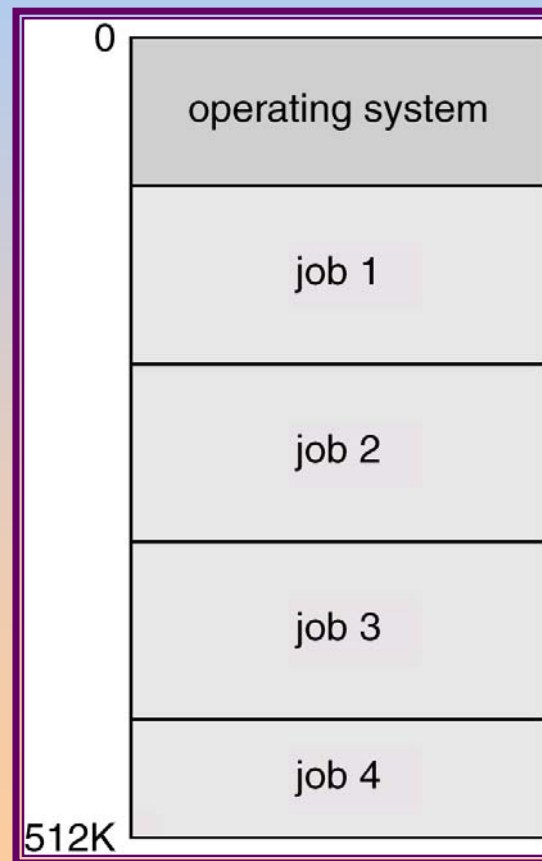  - when job completes control transfers pack to monitor

# Memory Layout for a Simple Batch System

```
┌─────────────────────┐
│                     │
│     operating       │
│      system         │
│                     │
├─────────────────────┤
│                     │
│                     │
│                     │
│                     │
│    user program     │
│       area          │
│                     │
│                     │
│                     │
│                     │
└─────────────────────┘
```

# Multiprogrammed Batch Systems

Several jobs are kept in main memory at the same time, and the CPU is multiplexed among them.

```
0
    operating system

    job 1

    job 2

    job 3

    job 4
512K
```

# OS Features Needed for Multiprogramming

- ■ I/O routine supplied by the system.
- ■ Memory management – the system must allocate the memory to several jobs.
- ■ CPU scheduling – the system must choose among several jobs ready to run.
- ■ Allocation of devices.

# Time-Sharing Systems–Interactive Computing

■ The CPU is multiplexed among several jobs that are kept in memory and on disk (the CPU is allocated to a job only if the job is in memory).

■ A job swapped in and out of memory to the disk.

■ On-line communication between the user and the system is provided; when the operating system finishes the execution of one command, it seeks the next "control statement" from the user's keyboard.

■ On-line system must be available for users to access data and code.

# Desktop Systems

- *Personal computers* – computer system dedicated to a single user.

- I/O devices – keyboards, mice, display screens, small printers.

- User convenience and responsiveness.

- Can adopt technology developed for larger operating system' often individuals have sole use of computer and do not need advanced CPU utilization of protection features.

- May run several different types of operating systems (Windows, MacOS, UNIX, Linux)

# Parallel Systems

- Multiprocessor systems with more than on CPU in close communication.

- *Tightly coupled system* – processors share memory and a clock; communication usually takes place through the shared memory.

- Advantages of parallel system:
  - Increased *throughput*
  - Economical
  - Increased reliability
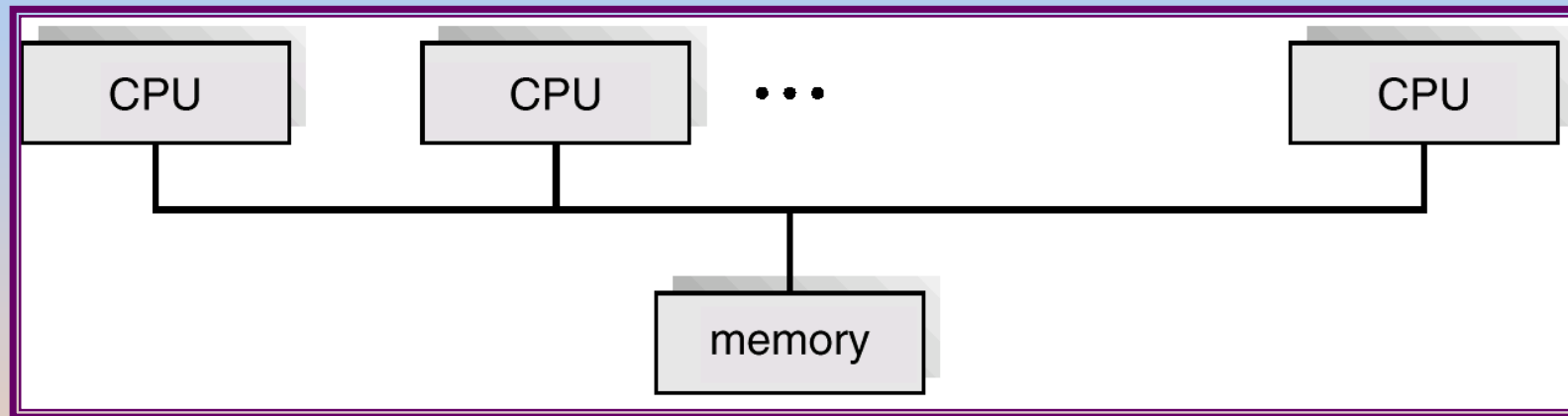    - ✔ graceful degradation
    - ✔ fail-soft systems

# Parallel Systems (Cont.)

- *Symmetric multiprocessing (SMP)*
    - Each processor runs and identical copy of the operating system.
    - Many processes can run at once without performance deterioration.
    - Most modern operating systems support SMP
- *Asymmetric multiprocessing*
    - Each processor is assigned a specific task; master processor schedules and allocated work to slave processors.
    - More common in extremely large systems

# Symmetric Multiprocessing Architecture

# Distributed Systems

- Distribute the computation among several physical processors.

- *Loosely coupled system* – each processor has its own local memory; processors communicate with one another through various communications lines, such as high-speed buses or telephone lines.

- Advantages of distributed systems.
  - Resources Sharing
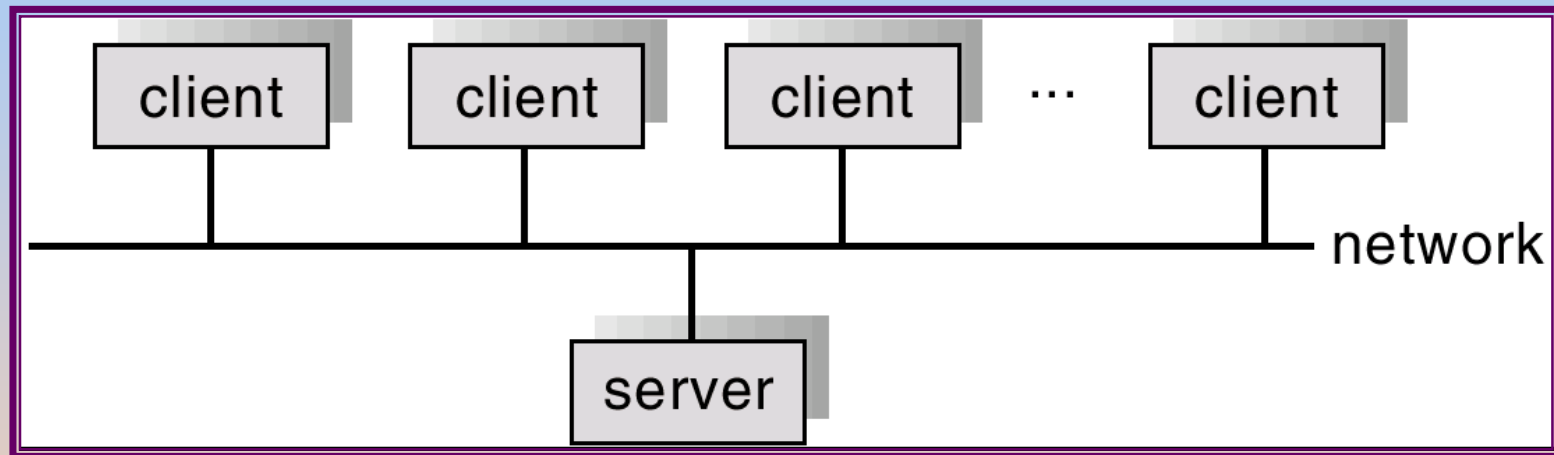  - Computation speed up – load sharing
  - Reliability
  - Communications

# Distributed Systems (cont)

- Requires networking infrastructure.
- Local area networks (LAN) or Wide area networks (WAN)
- May be either client-server or peer-to-peer systems.

# General Structure of Client-Server

# Clustered Systems

- Clustering allows two or more systems to share storage.
- Provides high reliability.
- *Asymmetric clustering*: one server runs the application while other servers standby.
- *Symmetric clustering*: all N hosts are running the application.

# Real-Time Systems

- Often used as a control device in a dedicated application such as controlling scientific experiments, medical imaging systems, industrial control systems, and some display systems.

- Well-defined fixed-time constraints.

- Real-Time systems may be either *hard* or *soft* real-time.

# Real-Time Systems (Cont.)

- **Hard real-time:**
  - Secondary storage limited or absent, data stored in short term memory, or read-only memory (ROM)
  - Conflicts with time-sharing systems, not supported by general-purpose operating systems.

- **Soft real-time**
  - Limited utility in industrial control of robotics
  - Useful in applications (multimedia, virtual reality) requiring advanced operating-system features.

# Handheld Systems

- Personal Digital Assistants (PDAs)
- Cellular telephones
- Issues:
    - Limited memory
    - Slow processors
    - Small display screens.

# Migration of Operating-System Concepts and Features

| | 1950 | 1960 | 1970 | 1980 | 1990 | 2000 |
|---|---|---|---|---|---|---|

mainframes — MULTICS

no software    compilers    time shared    multiuser

batch

resident monitors

networked

distributed systems

multiprocessor

fault tolerant

UNIX

minicomputers —

no software    compilers

time shared    multiuser    multiprocessor

resident monitors    networked    fault tolerant

clustered

UNIX

desktop computers —

no software    compilers    interactive    multiprocessor

multiuser    networked

UNIX

handheld computers —

compilers    no software

interactive

networked

# Computing Environments

- Traditional computing
- Web-Based Computing
- Embedded Computing

# Chapter 2:  Computer-System Structures

- Computer System Operation
- I/O Structure
- Storage Structure
- Storage Hierarchy
- Hardware Protection
- General System Architecture

# Computer-System Architecture

# 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*.

# Common Functions of Interrupts

- Interrupt transfers control to the interrupt service routine generally, through the *interrupt vector*, which contains the addresses of all the service routines.

- Interrupt architecture must save the 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 Handling

- The operating system preserves the state of the CPU by storing registers and the program counter.

- Determines which type of interrupt has occurred:
  - *polling*
  - *vectored* interrupt system

- Separate segments of code determine what action should be taken for each type of interrupt

# Interrupt Time Line For a Single Process Doing Output

# I/O Structure

- After I/O starts, control returns to user program only upon I/O completion.
    - Wait instruction idles the CPU until the next interrupt
    - Wait loop (contention for memory access).
    - At most one I/O request is outstanding at a time, no simultaneous I/O processing.
- After I/O starts, control returns to user program without waiting for I/O completion.
    - *System call* – request to the operating system to allow user to wait for I/O completion.
    - *Device-status table* contains entry for each I/O device indicating its type, address, and state.
    - Operating system indexes into I/O device table to determine device status and to modify table entry to include interrupt.

# Two I/O Methods

Synchronous

Asynchronous



(a)

(b)

# Device-Status Table



device: card reader 1
status: idle

device: line printer 3
status: busy

device: disk unit 1
status: idle

device: disk unit 2
status: idle

device: disk unit 3
status: busy

request for
line printer
address: 38546
length: 1372

request for
disk unit 3

file: xxx
operation: read
address: 43046
length: 20000

request for
disk unit 3

file: yyy
operation: write
address: 03458
length: 500

# Direct Memory Access Structure

- Used for high-speed I/O devices able to transmit information at close to memory speeds.

- Device controller transfers blocks of data from buffer storage directly to main memory without CPU intervention.

- Only on interrupt is generated per block, rather than the one interrupt per byte.

# 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.

# Moving-Head Disk Mechanism

# 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.

# Storage-Device Hierarchy

registers

cache

main memory

electronic disk

magnetic disk

optical disk

magnetic tapes

# Caching

- Use of high-speed memory to hold recently-accessed data.

- Requires a *cache management* policy.

- Caching introduces another level in storage hierarchy. This requires data that is simultaneously stored in more than one level to be *consistent*.

# Migration of A From Disk to Register

# Hardware Protection

- Dual-Mode Operation
- I/O Protection
- Memory Protection
- CPU Protection

# Dual-Mode Operation

- Sharing system resources requires operating system to ensure that an incorrect program cannot cause other programs to execute incorrectly.

- Provide hardware support to differentiate between at least two modes of operations.

  1. *User mode* – execution done on behalf of a user.
  2. *Monitor mode* (also *kernel mode* or *system mode*) – execution done on behalf of operating system.

# Dual-Mode Operation (Cont.)

- *Mode bit* added to computer hardware to indicate the current mode: monitor (0) or user (1).
- When an interrupt or fault occurs hardware switches to monitor mode.

Interrupt/fault

monitor          user

set user mode

*Privileged instructions* can be issued only in monitor mode.

# I/O Protection

- All I/O instructions are privileged instructions.

- Must ensure that a user program could never gain control of the computer in monitor mode (I.e., a user program that, as part of its execution, stores a new address in the interrupt vector).
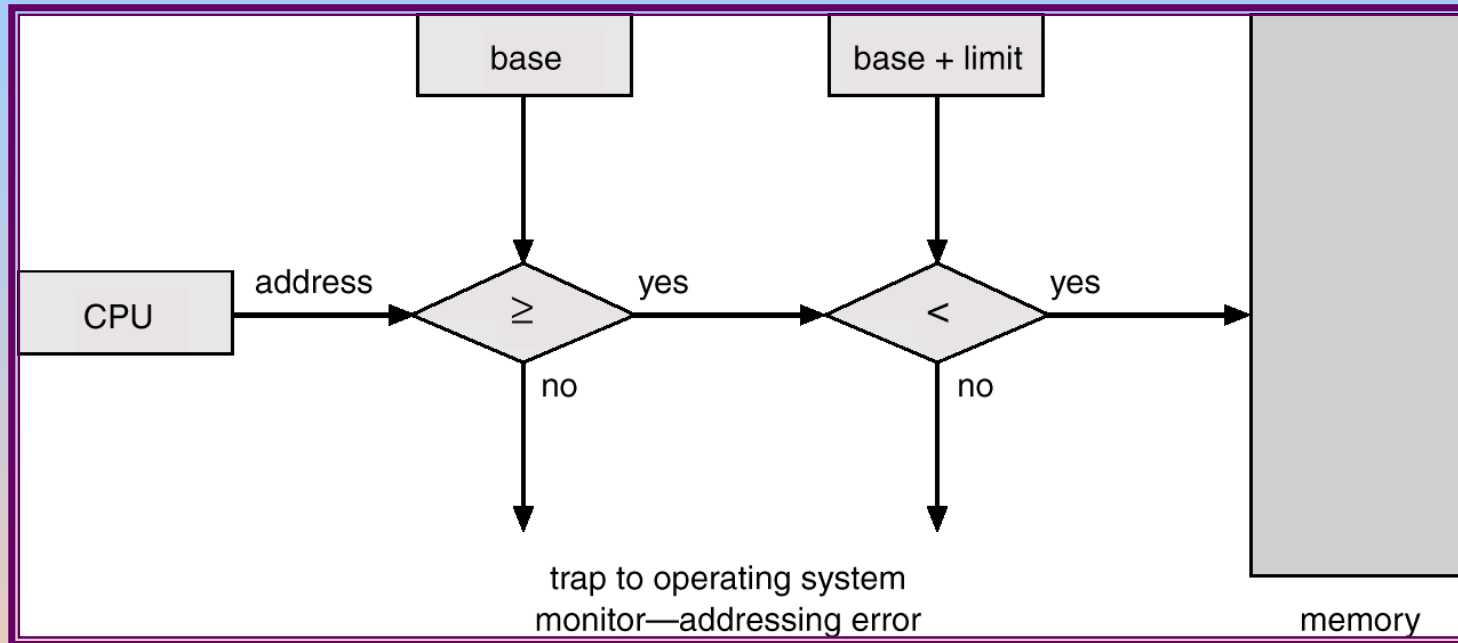
# Use of A System Call to Perform I/O

resident monitor

case *n*

① trap to monitor

② perform I/O

read

③ return to user

system call *n*

user program

# Memory Protection

- Must provide memory protection at least for the interrupt vector and the interrupt service routines.

- In order to have memory protection, add two registers that determine the range of legal addresses a program may access:

  - **Base register** – holds the smallest legal physical memory address.

  - **Limit register** – contains the size of the range

- Memory outside the defined range is protected.

# Use of A Base and Limit Register

# Hardware Address Protection

# Hardware Protection

- When executing in monitor mode, the operating system has unrestricted access to both monitor and user's memory.

- The load instructions for the *base* and *limit* registers are privileged instructions.
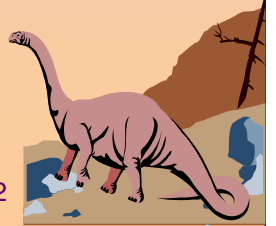
# CPU Protection

- *Timer* – interrupts computer after specified period to ensure operating system maintains control.
    - Timer is decremented every clock tick.
    - When timer reaches the value 0, an interrupt occurs.
- Timer commonly used to implement time sharing.
- Time also used to compute the current time.
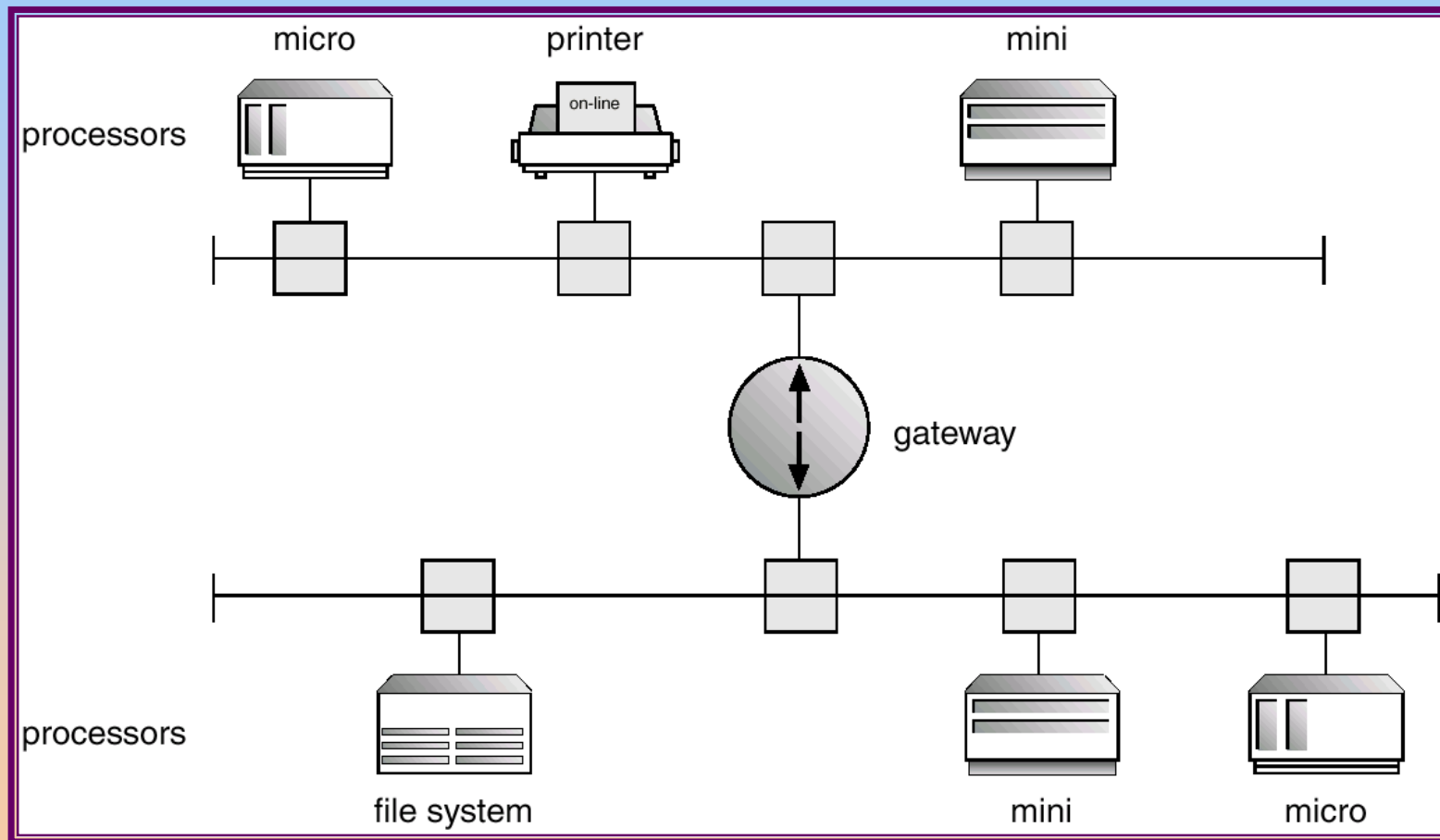- Load-timer is a privileged instruction.
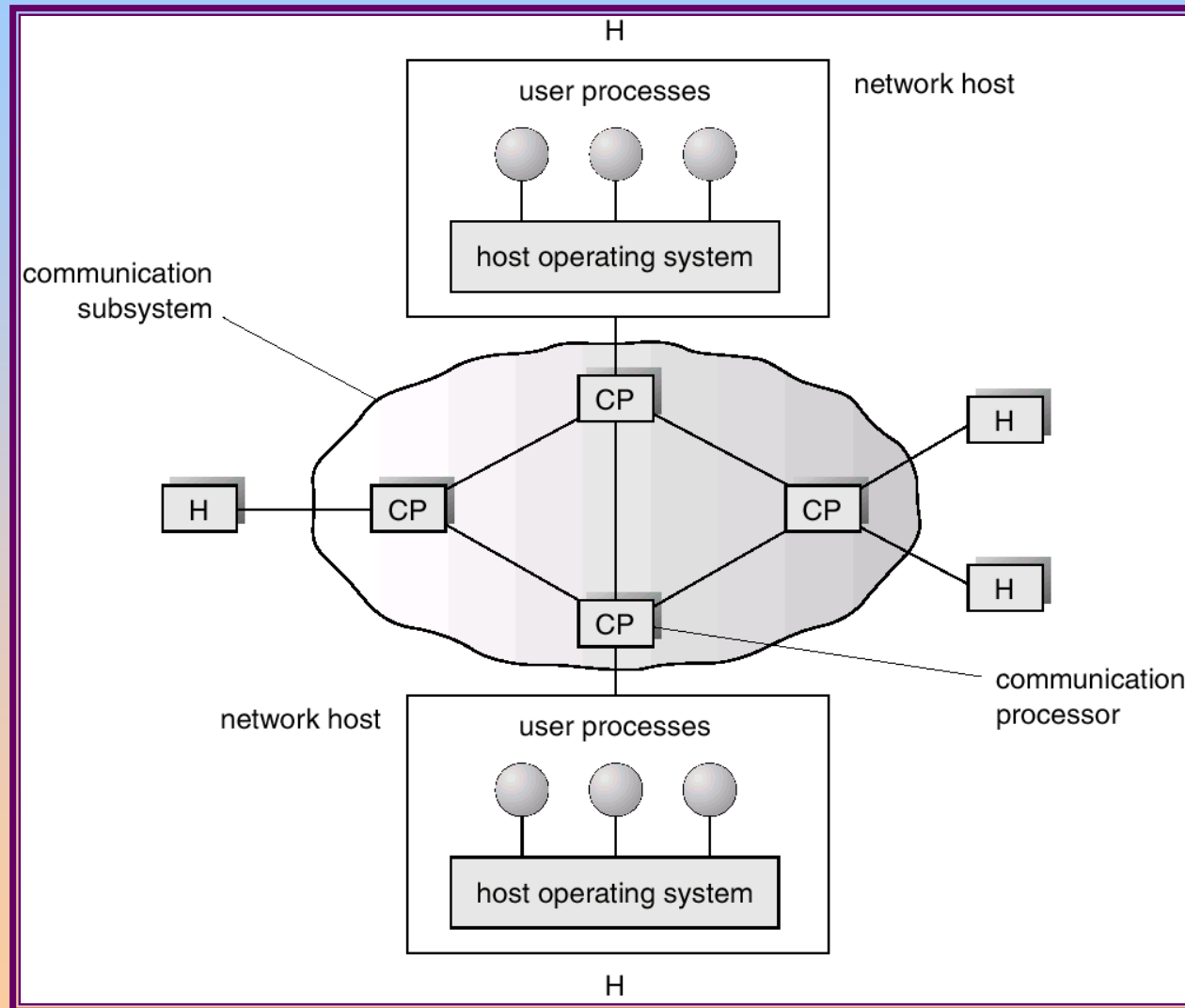
# Network Structure

- Local Area Networks (LAN)
- Wide Area Networks (WAN)

# Local Area Network Structure

# Wide Area Network Structure

# Chapter 3:  Operating-System Structures

- System Components
- Operating System Services
- System Calls
- System Programs
- System Structure
- Virtual Machines
- System Design and Implementation
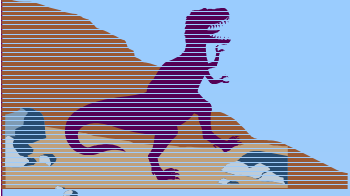- System Generation

# Common System Components

- Process Management

- Main Memory Management

- File Management

- I/O System Management

- Secondary Management

- Networking

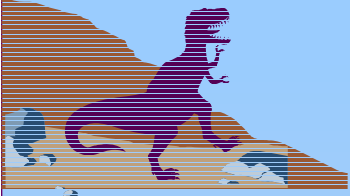- Protection System

- Command-Interpreter System

# Process Management

- A *process* is a program in execution.  A process needs certain resources, including CPU time, memory, files, and I/O devices, to accomplish its task.

- The operating system is responsible for the following activities in connection with process management.
  - Process creation and deletion.
  - process suspension and resumption.
  - Provision of mechanisms for:
    - process synchronization
    - process communication

# Main-Memory Management

- Memory is a large array of words or bytes, each with its own address. It is a repository of quickly accessible data shared by the CPU and I/O devices.

- Main memory is a volatile storage device. It loses its contents in the case of system failure.

- The operating system is responsible for the following activities in connections with memory management:

    - Keep track of which parts of memory are currently being used and by whom.

    - Decide which processes to load when memory space becomes available.

    - Allocate and deallocate memory space as needed.

# File Management

- A file is a collection of related information defined by its creator.  Commonly, files represent programs (both source and object forms) and data.

- The operating system is responsible for the following activities in connections with file management:

  - File creation and deletion.

  - Directory creation and deletion.

  - Support of primitives for manipulating files and directories.

  - Mapping files onto secondary storage.

  - File backup on stable (nonvolatile) storage media.

# I/O System Management

- The I/O system consists of:
    - A buffer-caching system
    - A general device-driver interface
    - Drivers for specific hardware devices

# Secondary-Storage Management

- Since main memory (*primary storage*) is volatile and too small to accommodate all data and programs permanently, the computer system must provide *secondary storage* to back up main memory.

- Most modern computer systems use disks as the principle on-line storage medium, for both programs and data.

- The operating system is responsible for the following activities in connection with disk management:
  - Free space management
  - Storage allocation
  - Disk scheduling

# Networking (Distributed Systems)

- A *distributed* system is a collection processors that do not share memory or a clock. Each processor has its own local memory.

- The processors in the system are connected through a communication network.

- Communication takes place using a *protocol.*

- A distributed system provides user access to various system resources.

- Access to a shared resource allows:
  - Computation speed-up
  - Increased data availability
  - Enhanced reliability

# Protection System

- *Protection* refers to a mechanism for controlling access by programs, processes, or users to both system and user resources.

- The protection mechanism must:
  - distinguish between authorized and unauthorized usage.
  - specify the controls to be imposed.
  - provide a means of enforcement.

# Command-Interpreter System

■ Many commands are given to the operating system by control statements which deal with:

  ✦ process creation and management

  ✦ I/O handling

  ✦ secondary-storage management

  ✦ main-memory management

  ✦ file-system access

  ✦ protection

  ✦ networking

# Command-Interpreter System (Cont.)

■ The program that reads and interprets control statements is called variously:

- ✦ command-line interpreter
- ✦ shell (in UNIX)

  Its function is to get and execute the next command statement.

# Operating System Services

- Program execution – system capability to load a program into memory and to run it.

- I/O operations – since user programs cannot execute I/O operations directly, the operating system must provide some means to perform I/O.

- File-system manipulation – program capability to read, write, create, and delete files.

- Communications – exchange of information between processes executing either on the same computer or on different systems tied together by a network. Implemented via *shared memory* or *message passing*.

- Error detection – ensure correct computing by detecting errors in the CPU and memory hardware, in I/O devices, or in user programs.

# Additional Operating System Functions

Additional functions exist not for helping the user, but rather for ensuring efficient system operations.
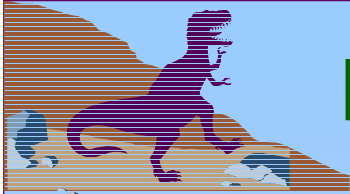
- Resource allocation – allocating resources to multiple users or multiple jobs running at the same time.

- Accounting – keep track of and record which users use how much and what kinds of computer resources for account billing or for accumulating usage statistics.
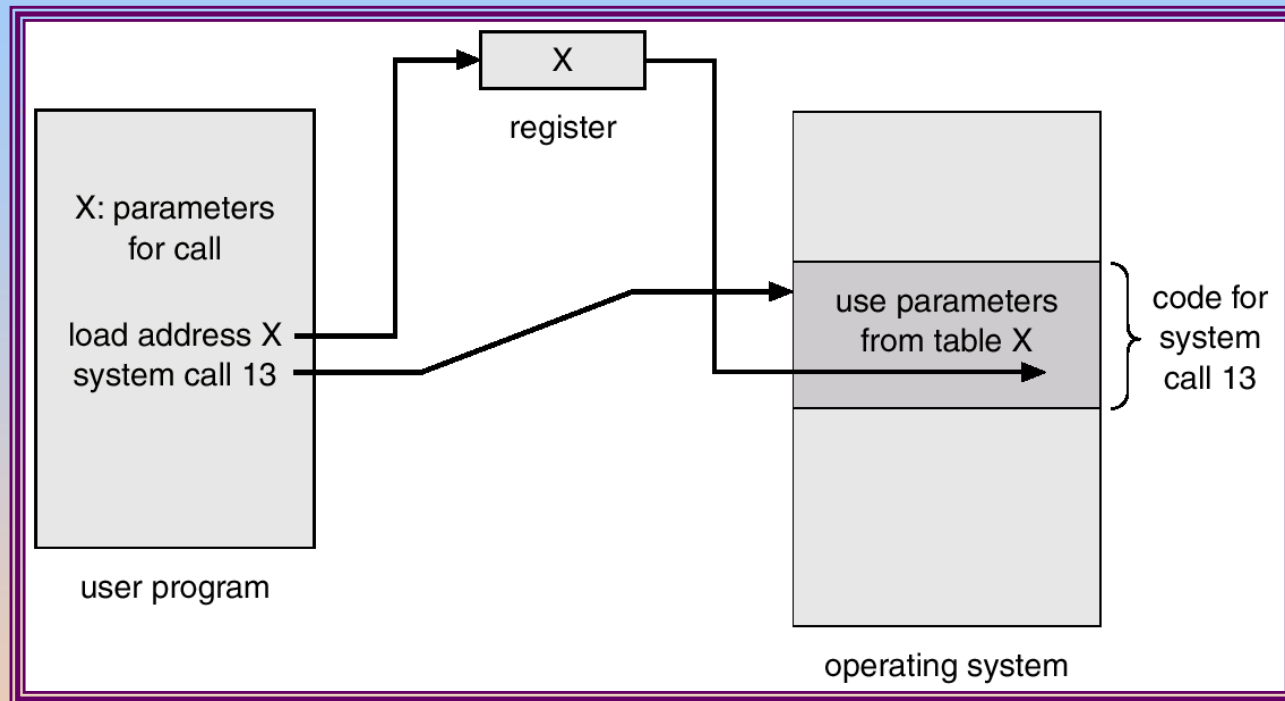
- Protection – ensuring that all access to system resources is controlled.

# System Calls

- System calls provide the interface between a running program and the operating system.
    - Generally available as assembly-language instructions.
    - Languages defined to replace assembly language for systems programming allow system calls to be made directly (e.g., C, C++)
- Three general methods are used to pass parameters between a running program and the operating system.
    - Pass parameters in *registers*.
    - Store the parameters in a table in memory, and the table address is passed as a parameter in a register.
    - *Push* (store) the parameters onto the *stack* by the program, and *pop* off the stack by operating system.
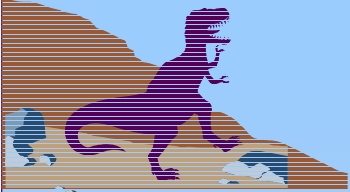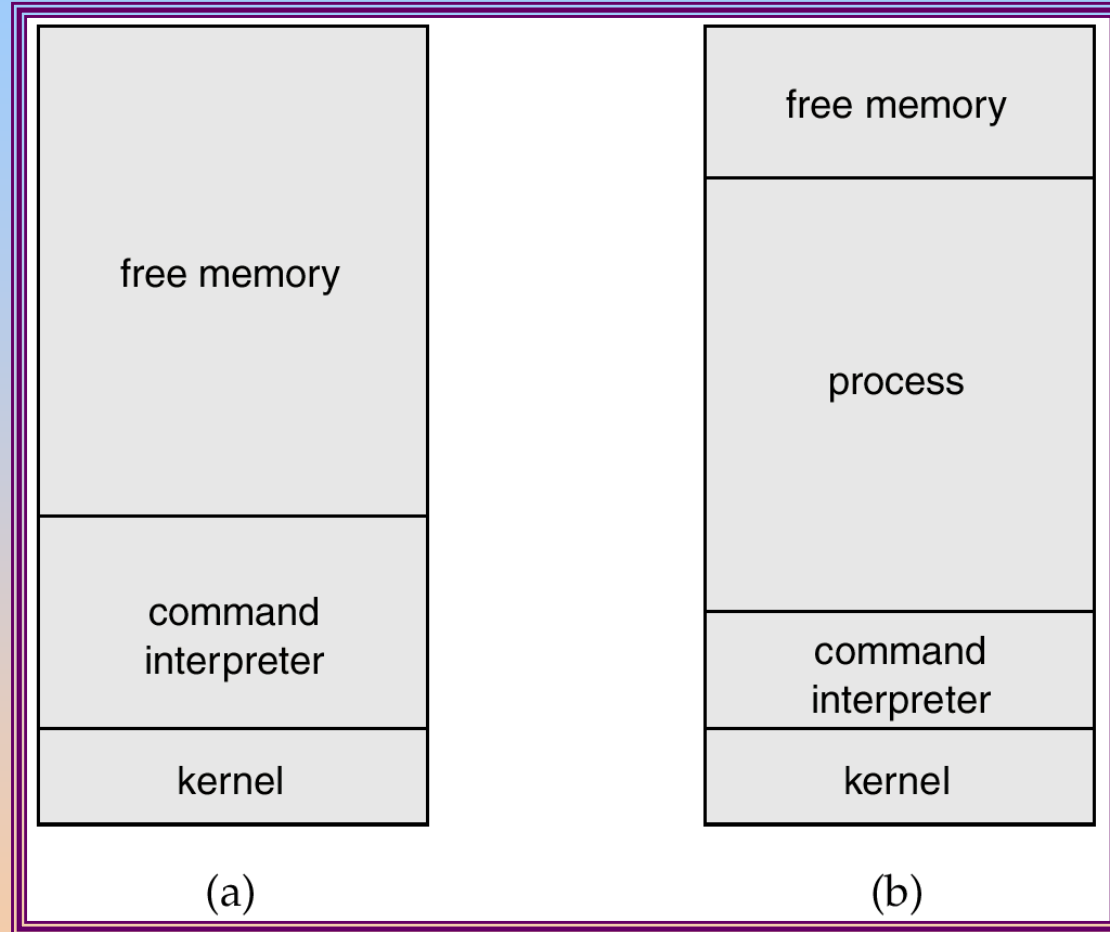
# Passing of Parameters As A Table

# Types of System Calls

- Process control
- File management
- Device management
- Information maintenance
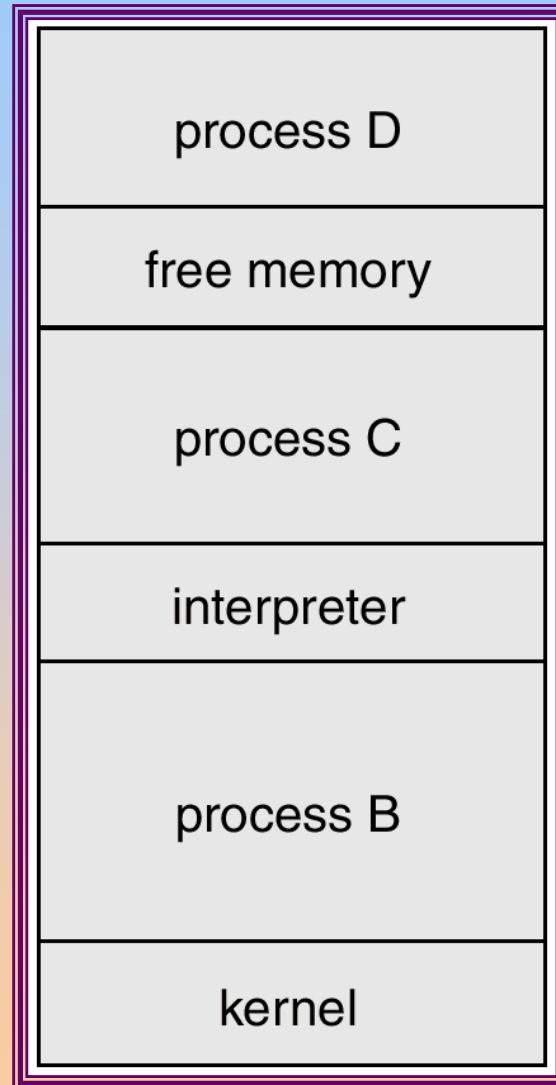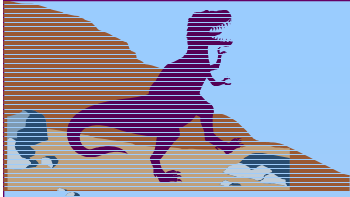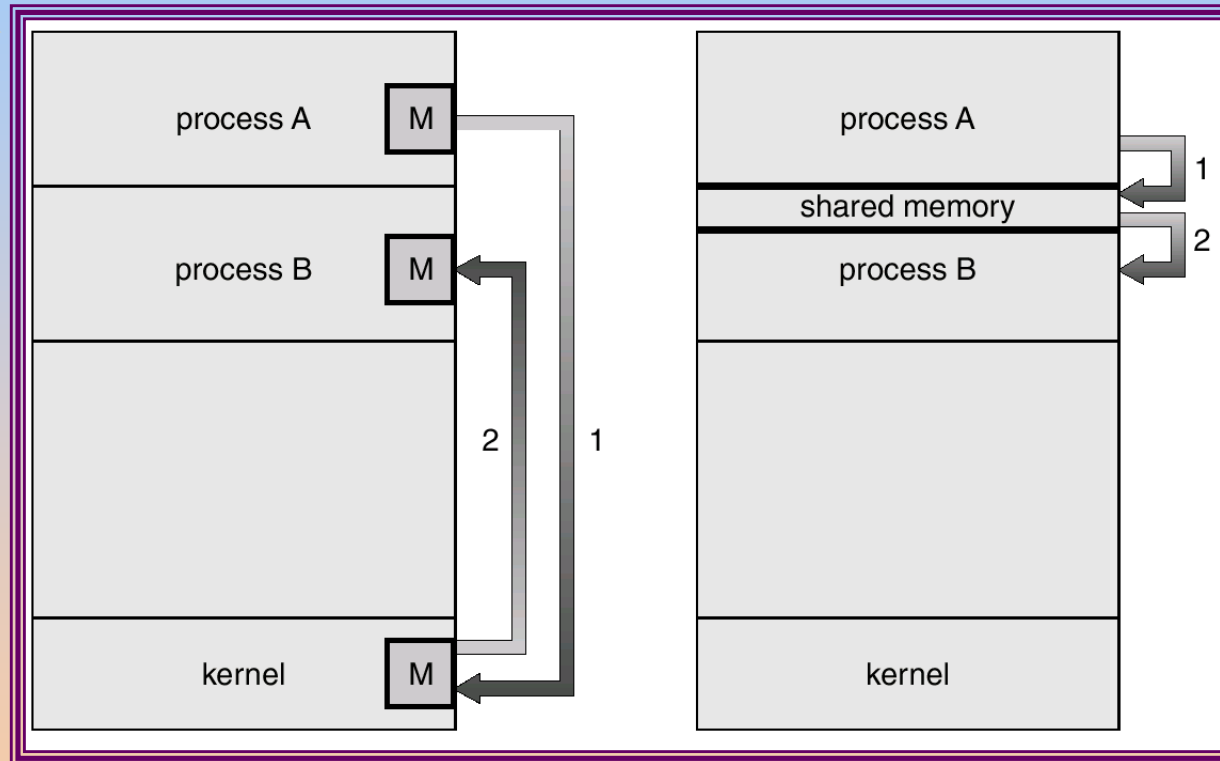- Communications

# MS-DOS Execution

| free memory |
|:---:|
| command interpreter |
| kernel |

(a)

At System Start-up

| free memory |
|:---:|
| process |
| command interpreter |
| kernel |

(b)

Running a Program

# UNIX Running Multiple Programs

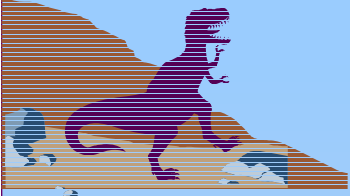| |
|---|
| process D |
| free memory |
| process C |
| interpreter |
| process B |
| kernel |

# Communication Models

- Communication may take place using either message passing or shared memory.



Msg Passing                    Shared Memory

# System Programs

- System programs provide a convenient environment for program development and execution. The can be divided into:
    - File manipulation
    - Status information
    - File modification
    - Programming language support
    - Program loading and execution
    - Communications
    - Application programs
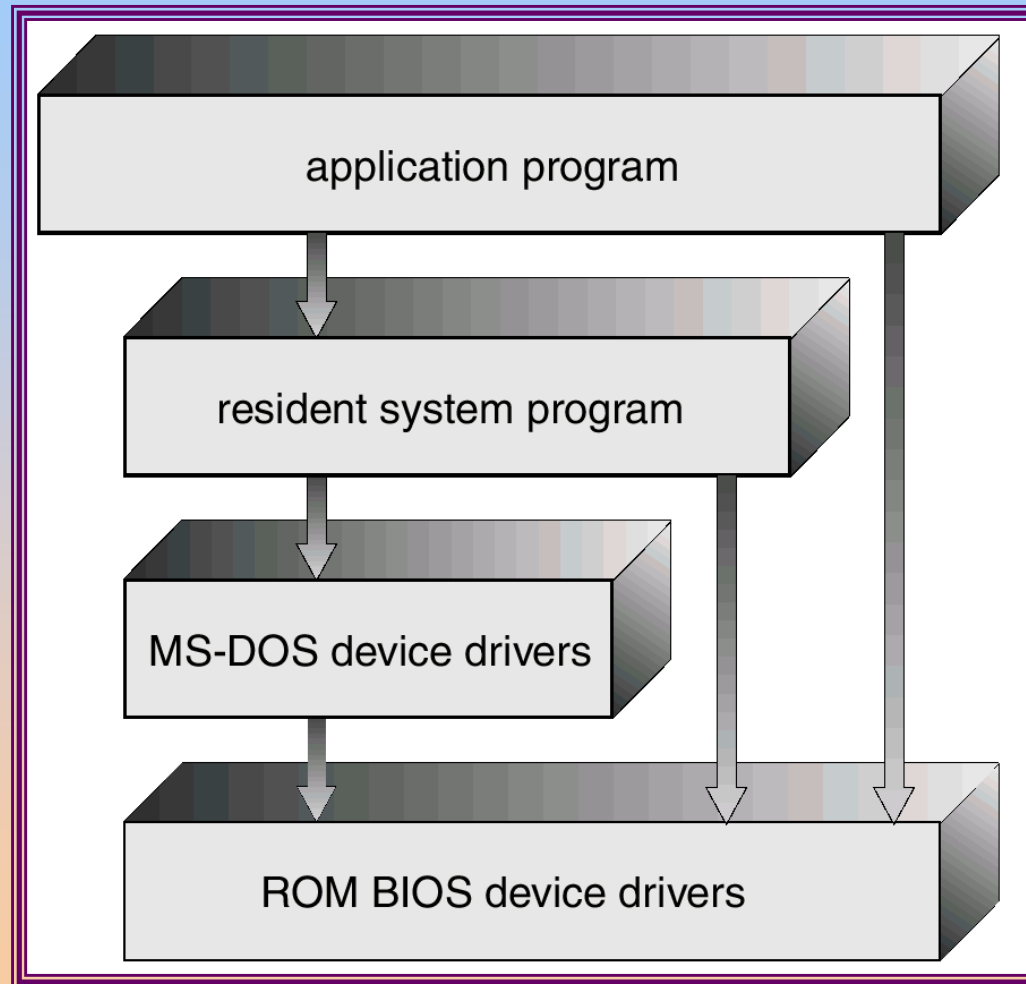- Most users' view of the operation system is defined by system programs, not the actual system calls.

# MS-DOS System Structure

- MS-DOS – written to provide the most functionality in the least space
  - not divided into modules
  - Although MS-DOS has some structure, its interfaces and levels of functionality are not well separated
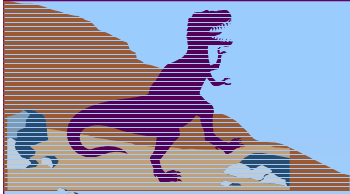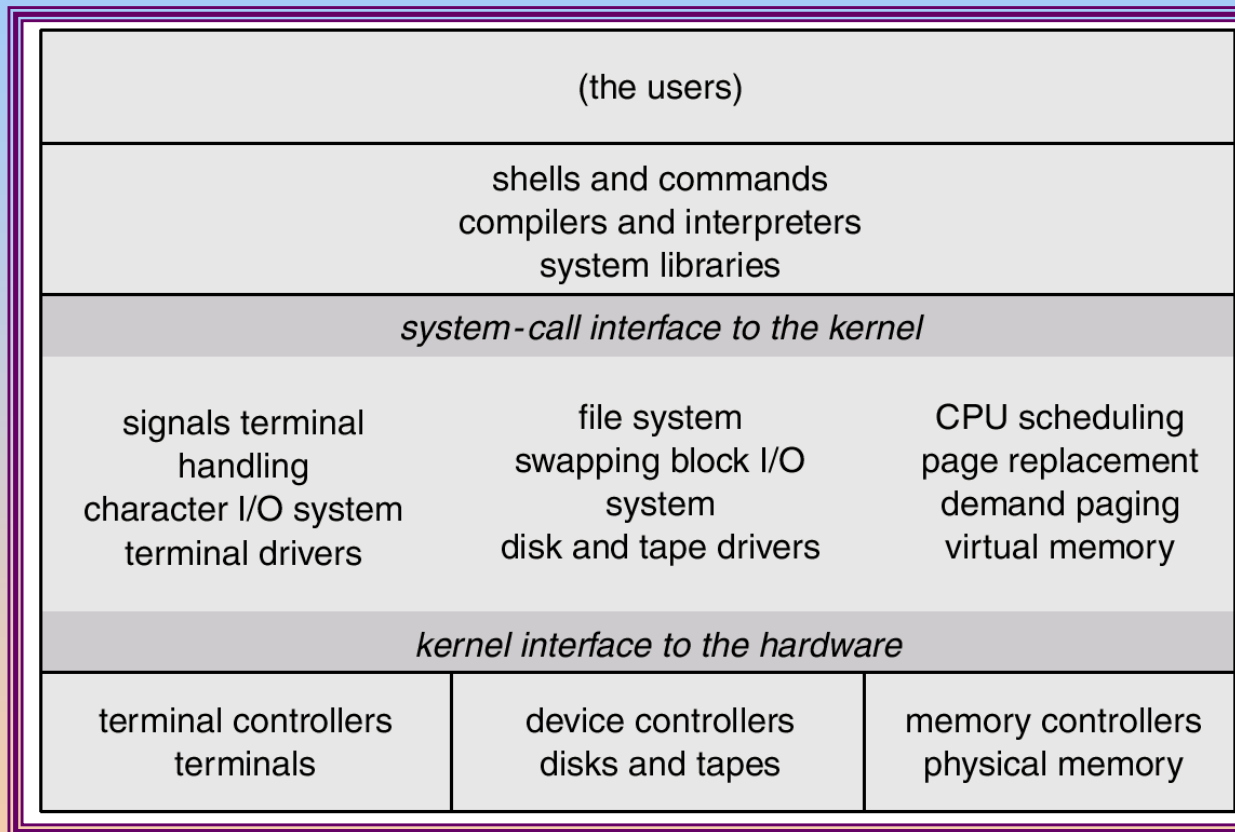
# MS-DOS Layer Structure

# UNIX System Structure

■ UNIX – limited by hardware functionality, the original UNIX operating system had limited structuring.  The UNIX OS consists of two separable parts.

✦ Systems programs

✦ The kernel

✔ Consists of everything below the system-call interface and above the physical hardware

✔ Provides the file system, CPU scheduling, memory management, and other operating-system functions; a large number of functions for one level.
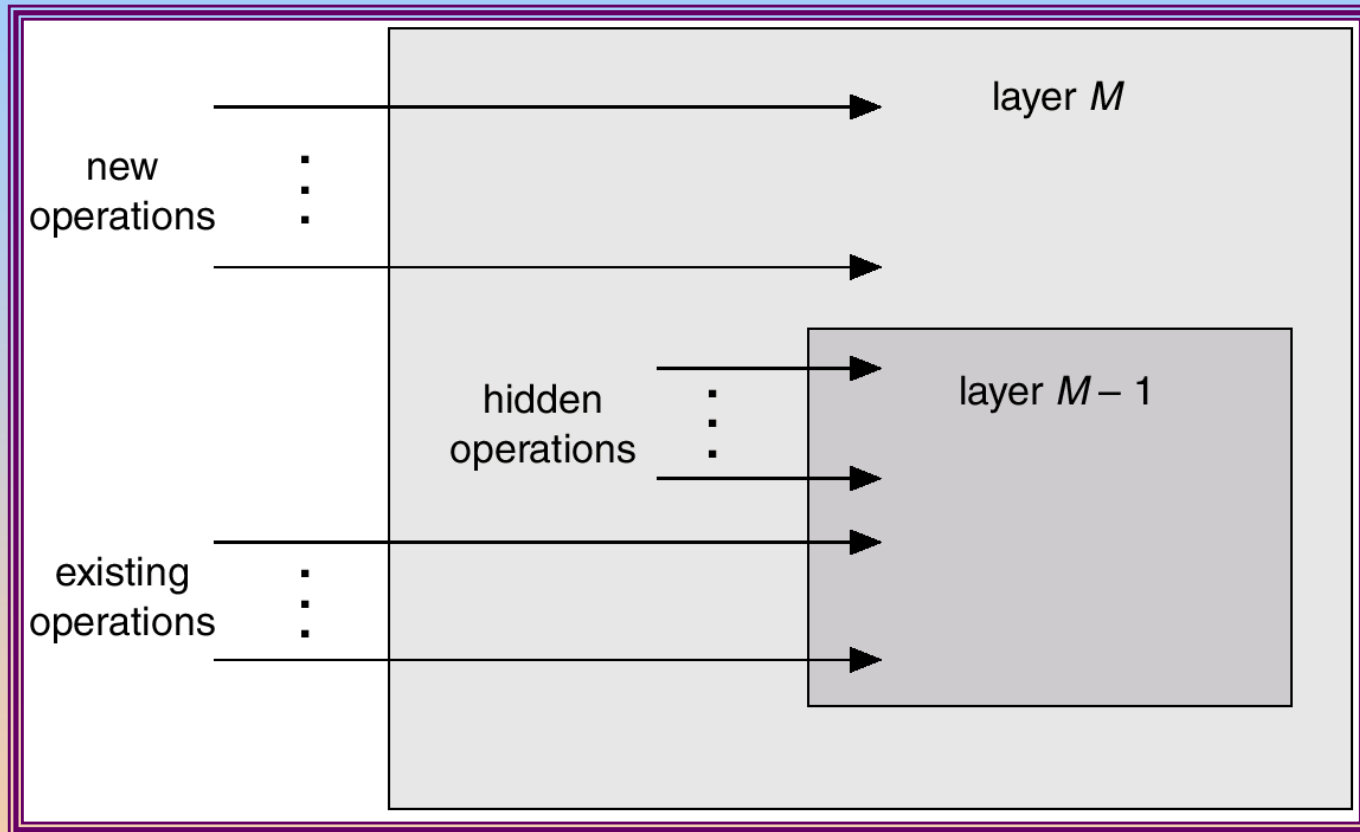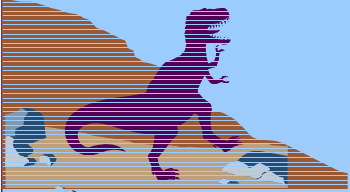
# UNIX System Structure

| (the users) |
| --- |
| shells and commands<br>compilers and interpreters<br>system libraries |
| *system-call interface to the kernel* |

| signals terminal<br>handling<br>character I/O system<br>terminal drivers | file system<br>swapping block I/O<br>system<br>disk and tape drivers | CPU scheduling<br>page replacement<br>demand paging<br>virtual memory |
| --- | --- | --- |

| *kernel interface to the hardware* | | |
| --- | --- | --- |
| terminal controllers<br>terminals | device controllers<br>disks and tapes | memory controllers<br>physical memory |

# Layered Approach

■ The operating system is divided into a number of layers (levels), each built on top of lower layers.  The bottom layer (layer 0), is the hardware; the highest (layer N) is the user interface.

■ With modularity, layers are selected such that each uses functions (operations) and services of only lower-level layers.
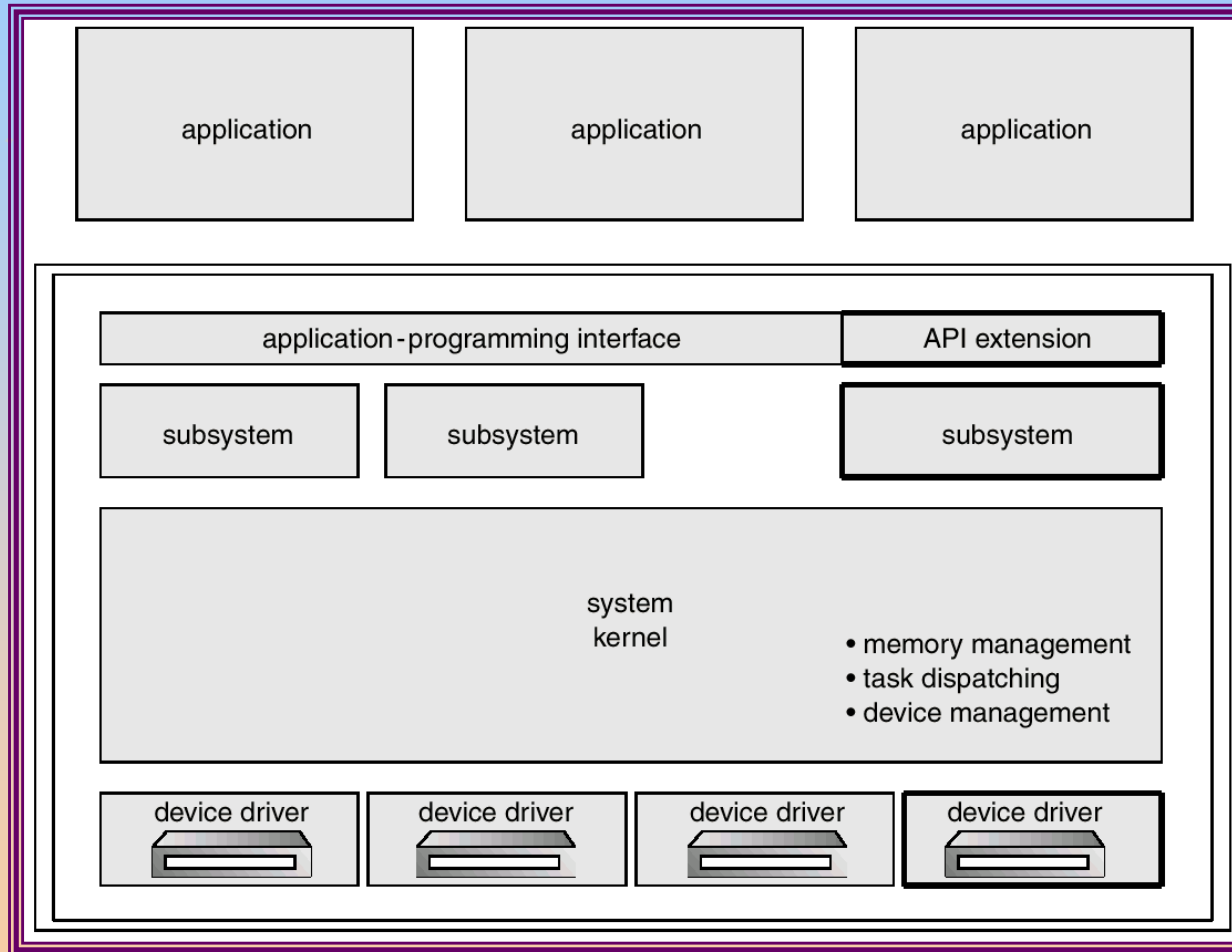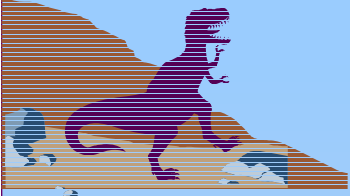
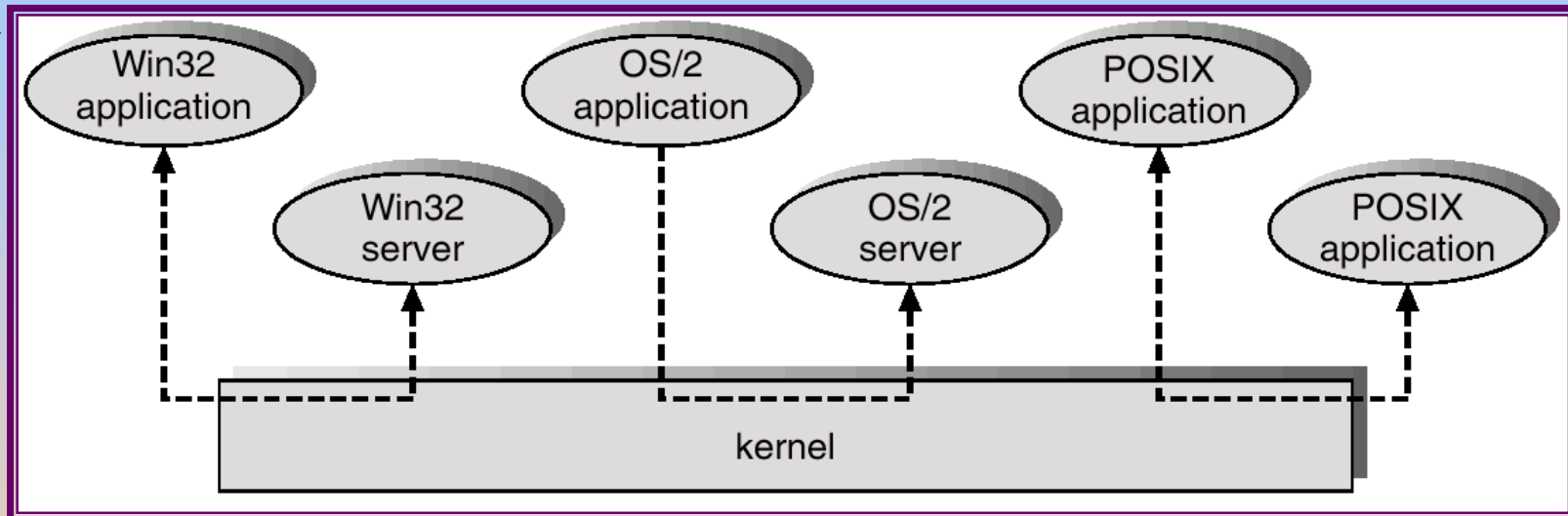# An Operating System Layer

# OS/2 Layer Structure

| application | application | application |
|:---:|:---:|:---:|

| application-programming interface | API extension |
|:---:|:---:|

| subsystem | subsystem | subsystem |
|:---:|:---:|:---:|

system
kernel

- memory management
- task dispatching
- device management

| device driver | device driver | device driver | device driver |
|:---:|:---:|:---:|:---:|

# Microkernel System Structure

- Moves as much from the kernel into "*user*" space.
- Communication takes place between user modules using message passing.
- Benefits:

  - easier to extend a microkernel

  - easier to port the operating system to new architectures

  - more reliable (less code is running in kernel mode)

  - more secure

# Windows NT Client-Server Structure

# Virtual Machines

- A *virtual machine* takes the layered approach to its logical conclusion. It treats hardware and the operating system kernel as though they were all hardware.

- A virtual machine provides an interface *identical* to the underlying bare hardware.

- The operating system creates the illusion of multiple processes, each executing on its own processor with its own (virtual) memory.
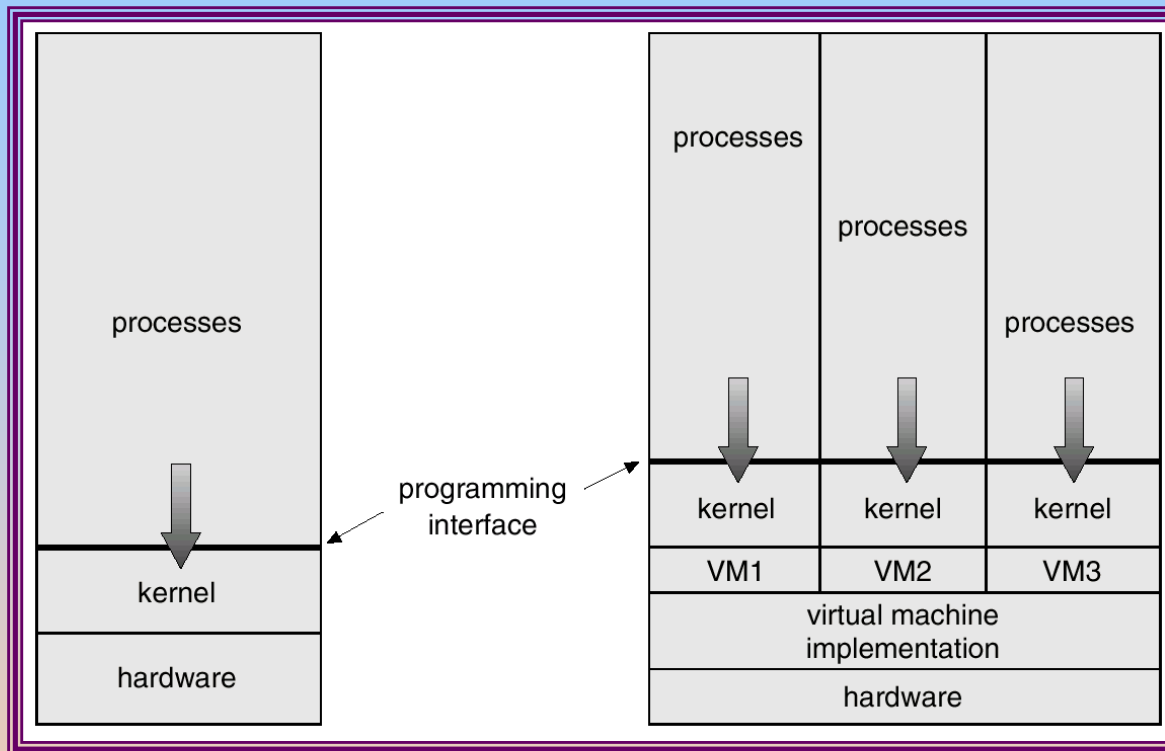
# Virtual Machines (Cont.)

- The resources of the physical computer are shared to create the virtual machines.

  - CPU scheduling can create the appearance that users have their own processor.

  - Spooling and a file system can provide virtual card readers and virtual line printers.

  - A normal user time-sharing terminal serves as the virtual machine operator's console.
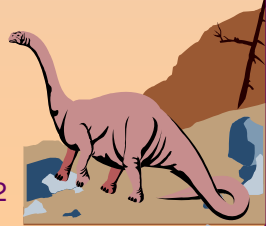
# System Models



Non-virtual Machine

Virtual Machine

# Advantages/Disadvantages of Virtual Machines

- The virtual-machine concept provides complete protection of system resources since each virtual machine is isolated from all other virtual machines.  This isolation, however, permits no direct sharing of resources.

- A virtual-machine system is a perfect vehicle for operating-systems research and development.  System development is done on the virtual machine, instead of on a physical machine and so does not disrupt normal system operation.

- The virtual machine concept is difficult to implement due to the effort required to provide an *exact* duplicate to the underlying machine.
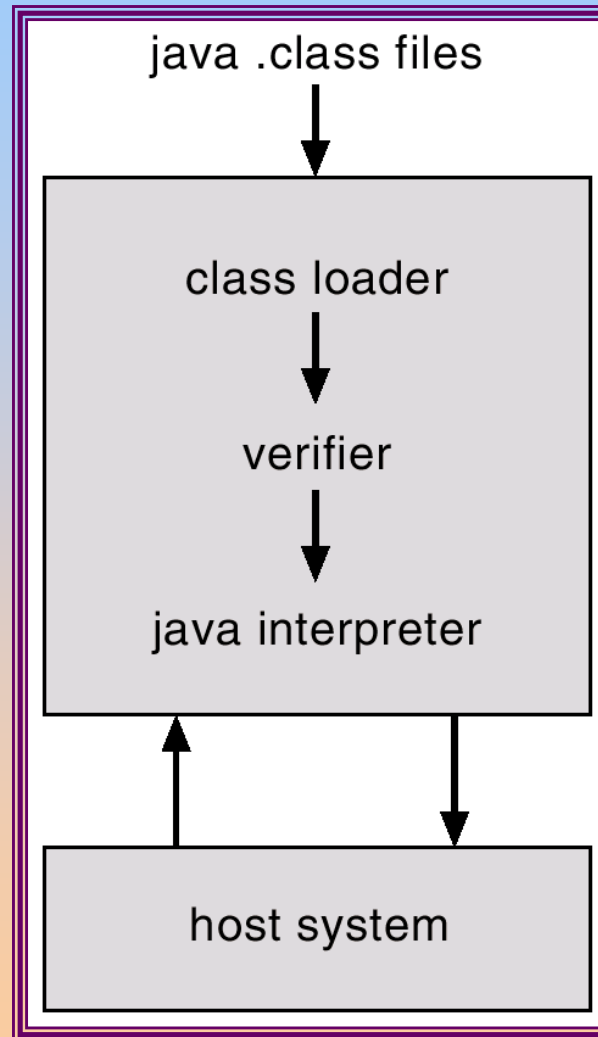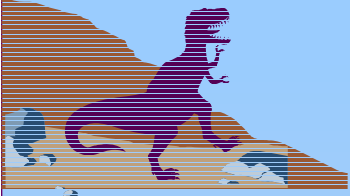
# Java Virtual Machine

- Compiled Java programs are platform-neutral bytecodes executed by a Java Virtual Machine (JVM).

- JVM consists of

  - class loader

  - class verifier

  - runtime interpreter

- Just-In-Time (JIT) compilers increase performance

# Java Virtual Machine

```
java .class files
        |
        v
+---------------------------+
|                           |
|       class loader        |
|           |               |
|           v               |
|        verifier           |
|           |               |
|           v               |
|      java interpreter     |
|                           |
+---------------------------+
      ^             |
      |             v
+---------------------------+
|       host system         |
+---------------------------+
```

# System Design Goals

■ User goals – operating system should be convenient to use, easy to learn, reliable, safe, and fast.

■ System goals – operating system should be easy to design, implement, and maintain, as well as flexible, reliable, error-free, and efficient.

# Mechanisms and Policies

- Mechanisms determine how to do something, policies decide what will be done.

- The separation of policy from mechanism is a very important principle, it allows maximum flexibility if policy decisions are to be changed later.

# System Implementation

- Traditionally written in assembly language, operating systems can now be written in higher-level languages.

- Code written in a high-level language:
  - can be written faster.
  - is more compact.
  - is easier to understand and debug.

- An operating system is far easier to *port* (move to some other hardware) if it is written in a high-level language.

# System Generation (SYSGEN)

- Operating systems are designed to run on any of a class of machines; the system must be configured for each specific computer site.

- SYSGEN program obtains information concerning the specific configuration of the hardware system.

- *Booting* – starting a computer by loading the kernel.

- *Bootstrap program* – code stored in ROM that is able to locate the kernel, load it into memory, and start its execution.

# Chapter 4:  Processes

- Process Concept
- Process Scheduling
- Operations on Processes
- Cooperating Processes
- Interprocess Communication
- Communication in Client-Server Systems

# Process Concept

- An operating system executes a variety of programs:
    - Batch system – jobs
    - Time-shared systems – user programs or tasks
- Textbook uses the terms *job* and *process* almost interchangeably.
- Process – a program in execution; process execution must progress in sequential fashion.
- A process includes:
    - program counter
    - stack
    - data section

# Process State

- As a process executes, it changes *state*
  - **new**:  The process is being created.
  - **running**:  Instructions are being executed.
  - **waiting**:  The process is waiting for some event to occur.
  - **ready**:  The process is waiting to be assigned to a process.
  - **terminated**:  The process has finished execution.

# Diagram of Process State

# Process Control Block (PCB)

Information associated with each process.

- Process state
- Program counter
- CPU registers
- CPU scheduling information
- Memory-management information
- Accounting information
- I/O status information

# Process Control Block (PCB)

| pointer | process state |
|---------|---------------|
| process number | |
| program counter | |
| registers | |
| memory limits | |
| list of open files | |
| ⋮ | |

# CPU Switch From Process to Process

| process $P_0$ | operating system | process $P_1$ |
|---|---|---|

interrupt or system call

executing

save state into $PCB_0$

⋮

reload state from $PCB_1$

idle

executing

interrupt or system call

idle

save state into $PCB_1$

⋮

reload state from $PCB_0$

idle

executing

# Process Scheduling Queues

- ■ Job queue – set of all processes in the system.

- ■ Ready queue – set of all processes residing in main memory, ready and waiting to execute.

- ■ Device queues – set of processes waiting for an I/O device.

- ■ Process migration between the various queues.

# Representation of Process Scheduling

# Schedulers

- Long-term scheduler (or job scheduler) – selects which processes should be brought into the ready queue.

- Short-term scheduler (or CPU scheduler) – selects which process should be executed next and allocates CPU.

# Addition of Medium Term Scheduling

# Schedulers (Cont.)

- Short-term scheduler is invoked very frequently (milliseconds) $\Rightarrow$ (must be fast).

- Long-term scheduler is invoked very infrequently (seconds, minutes) $\Rightarrow$ (may be slow).

- The long-term scheduler controls the *degree of multiprogramming.*

- Processes can be described as either:
  - I/O-*bound process* – spends more time doing I/O than computations, many short CPU bursts.
  - *CPU-bound process* – spends more time doing computations; few very long CPU bursts.

# Context Switch

- When CPU switches to another process, the system must save the state of the old process and load the saved state for the new process.

- Context-switch time is overhead; the system does no useful work while switching.

- Time dependent on hardware support.

# Process Creation

- Parent process create children processes, which, in turn create other processes, forming a tree of processes.

- Resource sharing
  - Parent and children share all resources.
  - Children share subset of parent's resources.
  - Parent and child share no resources.

- Execution
  - Parent and children execute concurrently.
  - Parent waits until children terminate.

# Process Creation (Cont.)

- Address space
  - ✦ Child duplicate of parent.
  - ✦ Child has a program loaded into it.
- UNIX examples
  - ✦ **fork** system call creates new process
  - ✦ **exec** system call used after a **fork** to replace the process' memory space with a new program.

# Processes Tree on a UNIX System

# Process Termination

- Process executes last statement and asks the operating system to decide it (**exit**).
    - ✦ Output data from child to parent (via **wait**).
    - ✦ Process' resources are deallocated by operating system.
- Parent may terminate execution of children processes (**abort**).
    - ✦ Child has exceeded allocated resources.
    - ✦ Task assigned to child is no longer required.
    - ✦ Parent is exiting.
        - ✔ Operating system does not allow child to continue if its parent terminates.
        - ✔ Cascading termination.

# Cooperating Processes

- *Independent* process cannot affect or be affected by the execution of another process.

- *Cooperating* process can affect or be affected by the execution of another process

- Advantages of process cooperation
    - Information sharing
    - Computation speed-up
    - Modularity
    - Convenience

# **Producer-Consumer Problem**

- Paradigm for cooperating processes, *producer* process produces information that is consumed by a *consumer* process.

    - *unbounded-buffer* places no practical limit on the size of the buffer.
    - *bounded-buffer* assumes that there is a fixed buffer size.

# Bounded-Buffer – Shared-Memory Solution

- Shared data

   ```
   #define BUFFER_SIZE 10
   Typedef struct {
      . . .
   } item;
   item buffer[BUFFER_SIZE];
   int in = 0;
   int out = 0;
   ```

- Solution is correct, but can only use BUFFER_SIZE-1 elements

# Bounded-Buffer – Producer Process

```
item nextProduced;

while (1) {
        while (((in + 1) % BUFFER_SIZE) == out)
                ; /* do nothing */
        buffer[in] = nextProduced;
        in = (in + 1) % BUFFER_SIZE;
}
```

# Bounded-Buffer – Consumer Process

```
item nextConsumed;

while (1) {
        while (in == out)
                ; /* do nothing */
        nextConsumed = buffer[out];
        out = (out + 1) % BUFFER_SIZE;
}
```

# Interprocess Communication (IPC)

- Mechanism for processes to communicate and to synchronize their actions.

- Message system – processes communicate with each other without resorting to shared variables.

- IPC facility provides two operations:
  - **send**(*message*) – message size fixed or variable
  - **receive**(*message*)

- If $P$ and $Q$ wish to communicate, they need to:
  - establish a *communication link* between them
  - exchange messages via send/receive

- Implementation of communication link
  - physical (e.g., shared memory, hardware bus)
  - logical (e.g., logical properties)

# Implementation Questions

- How are links established?

- Can a link be associated with more than two processes?

- How many links can there be between every pair of communicating processes?

- What is the capacity of a link?

- Is the size of a message that the link can accommodate fixed or variable?

- Is a link unidirectional or bi-directional?

# Direct Communication

- Processes must name each other explicitly:
    - **send** (*P, message*) – send a message to process P
    - **receive**(*Q, message*) – receive a message from process Q
- Properties of communication link
    - Links are established automatically.
    - A link is associated with exactly one pair of communicating processes.
    - Between each pair there exists exactly one link.
    - The link may be unidirectional, but is usually bi-directional.

# Indirect Communication

- Messages are directed and received from mailboxes (also referred to as ports).
    - Each mailbox has a unique id.
    - Processes can communicate only if they share a mailbox.
- Properties of communication link
    - Link established only if processes share a common mailbox
    - A link may be associated with many processes.
    - Each pair of processes may share several communication links.
    - Link may be unidirectional or bi-directional.

# Indirect Communication

- Operations
    - create a new mailbox
    - send and receive messages through mailbox
    - destroy a mailbox
- Primitives are defined as:

    **send**(*A, message*) – send a message to mailbox A

    **receive**(*A, message*) – receive a message from mailbox A

# Indirect Communication

- **Mailbox sharing**
  - $P_1$, $P_2$, and $P_3$ share mailbox A.
  - $P_1$, sends; $P_2$ and $P_3$ receive.
  - Who gets the message?

- **Solutions**
  - Allow a link to be associated with at most two processes.
  - Allow only one process at a time to execute a receive operation.
  - Allow the system to select arbitrarily the receiver.  Sender is notified who the receiver was.

# Synchronization

- Message passing may be either blocking or non-blocking.
- **Blocking** is considered **synchronous**
- **Non-blocking** is considered **asynchronous**
- **send** and **receive** primitives may be either blocking or non-blocking.

# Buffering

- Queue of messages attached to the link; implemented in one of three ways.

    1. Zero capacity – 0 messages
       Sender must wait for receiver (rendezvous).

    2. Bounded capacity – finite length of $n$ messages
       Sender must wait if link full.

    3. Unbounded capacity – infinite length
       Sender never waits.

# Client-Server Communication

- Sockets
- Remote Procedure Calls
- Remote Method Invocation (Java)

# Sockets

- A socket is defined as an *endpoint for communication*.

- Concatenation of IP address and port

- The socket **161.25.19.8:1625** refers to port **1625** on host **161.25.19.8**

- Communication consists between a pair of sockets.

# Socket Communication

host *X*

(146.86.5.20)

socket

(146.86.5.2/1625)

web server

(161.25.19.8)

socket

(161.25.19.8/80)

# Remote Procedure Calls

- Remote procedure call (RPC) abstracts procedure calls between processes on networked systems.
- **Stubs** – client-side proxy for the actual procedure on the server.
- The client-side stub locates the server and *marshalls* the parameters.
- The server-side stub receives this message, unpacks the marshalled parameters, and peforms the procedure on the server.

# Execution of RPC

# Remote Method Invocation

- Remote Method Invocation (RMI) is a Java mechanism similar to RPCs.

- RMI allows a Java program on one machine to invoke a method on a remote object.

# Marshalling Parameters

client

val = server.someMethod(A,B)

```
stub
```

remote object

boolean someMethod (Object x, Object y)
{
    implementation of someMethod
    . . .
}

```
skeleton
```

A, B, someMethod

boolean return value

4.38

# Chapter 5: Threads

- Overview
- Multithreading Models
- Threading Issues
- Pthreads
- Solaris 2 Threads
- Windows 2000 Threads
- Linux Threads
- Java Threads

# Single and Multithreaded Processes



| code | data | files |
| --- | --- | --- |
| registers | | stack |

thread →

single-threaded

| code | data | files |
| --- | --- | --- |
| registers | registers | registers |
| stack | stack | stack |

← thread

multithreaded

# Benefits

- Responsiveness

- Resource Sharing

- Economy

- Utilization of MP Architectures

# User Threads

■ Thread management done by user-level threads library

■ Examples
- POSIX *Pthreads*
- Mach *C-threads*
- Solaris *threads*

# Kernel Threads

■ Supported by the Kernel

■ Examples
- Windows 95/98/NT/2000
- Solaris
- Tru64 UNIX
- BeOS
- Linux

# Multithreading Models

- Many-to-One

- One-to-One

- Many-to-Many

# Many-to-One

- Many user-level threads mapped to single kernel thread.

- Used on systems that do not support kernel threads.

# Many-to-One Model

user thread

k   kernel thread

# One-to-One

■ Each user-level thread maps to kernel thread.

■ Examples
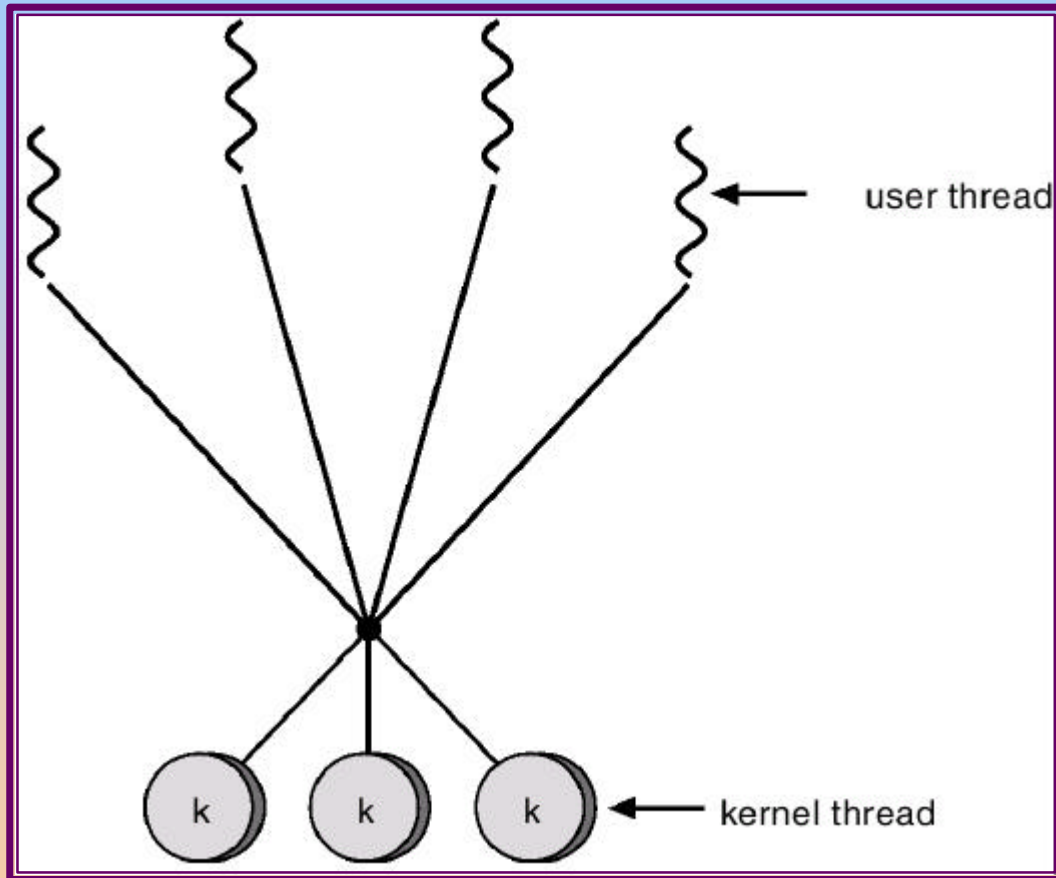  - Windows 95/98/NT/2000
  - OS/2

# One-to-one Model

# Many-to-Many Model

- Allows many user level threads to be mapped to many kernel threads.
- Allows the operating system to create a sufficient number of kernel threads.
- Solaris 2
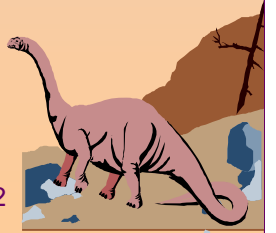- Windows NT/2000 with the *ThreadFiber* package

# Many-to-Many Model



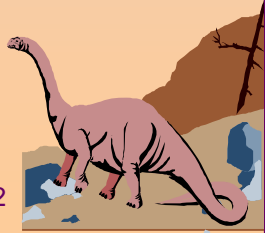user thread

kernel thread

# Threading Issues

- Semantics of fork() and exec() system calls.
- Thread cancellation.
- Signal handling
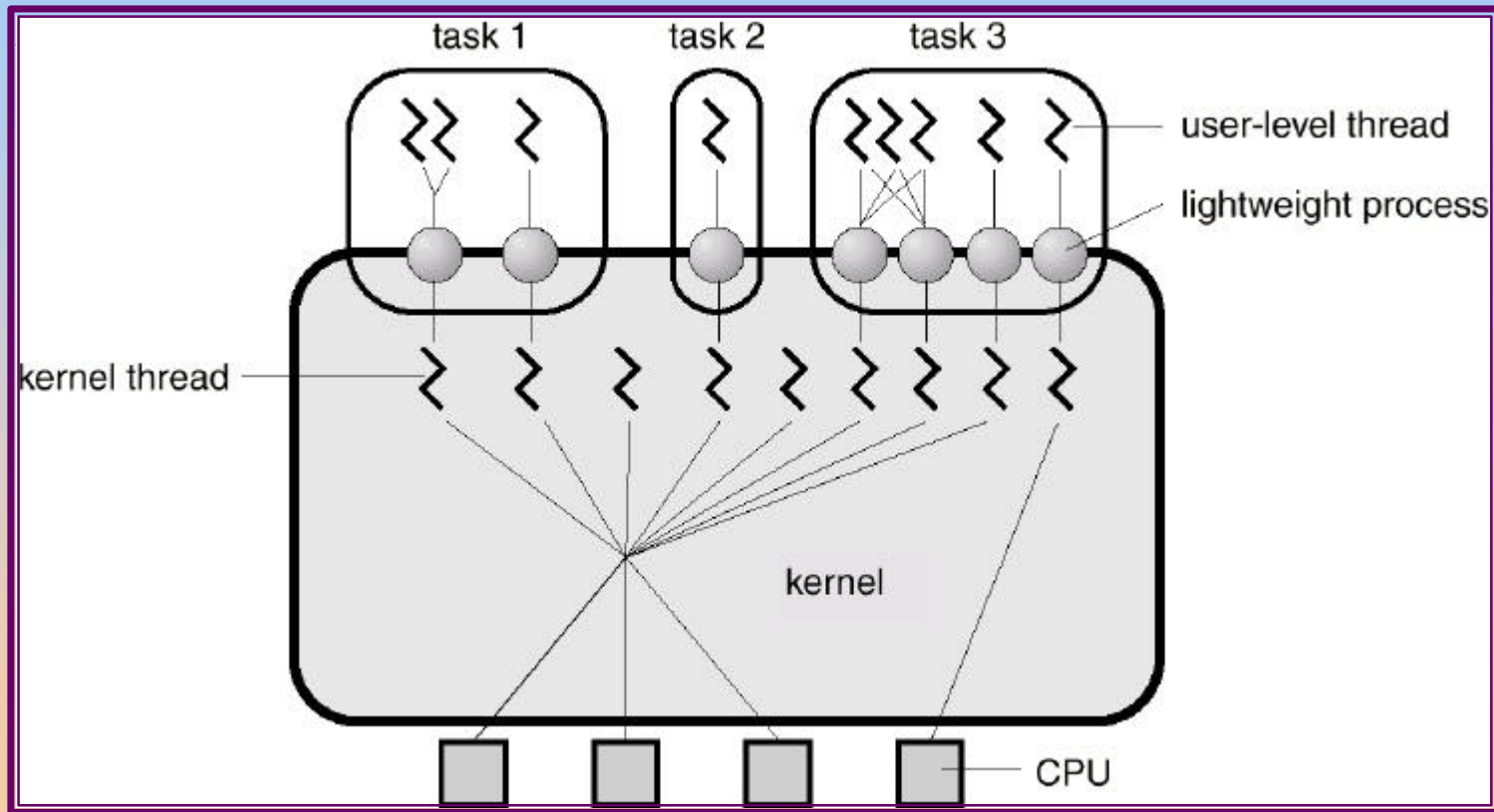- Thread pools
- Thread specific data

# Pthreads

- a POSIX standard (IEEE 1003.1c) API for thread creation and synchronization.

- API specifies behavior of the thread library, implementation is up to development of the library.
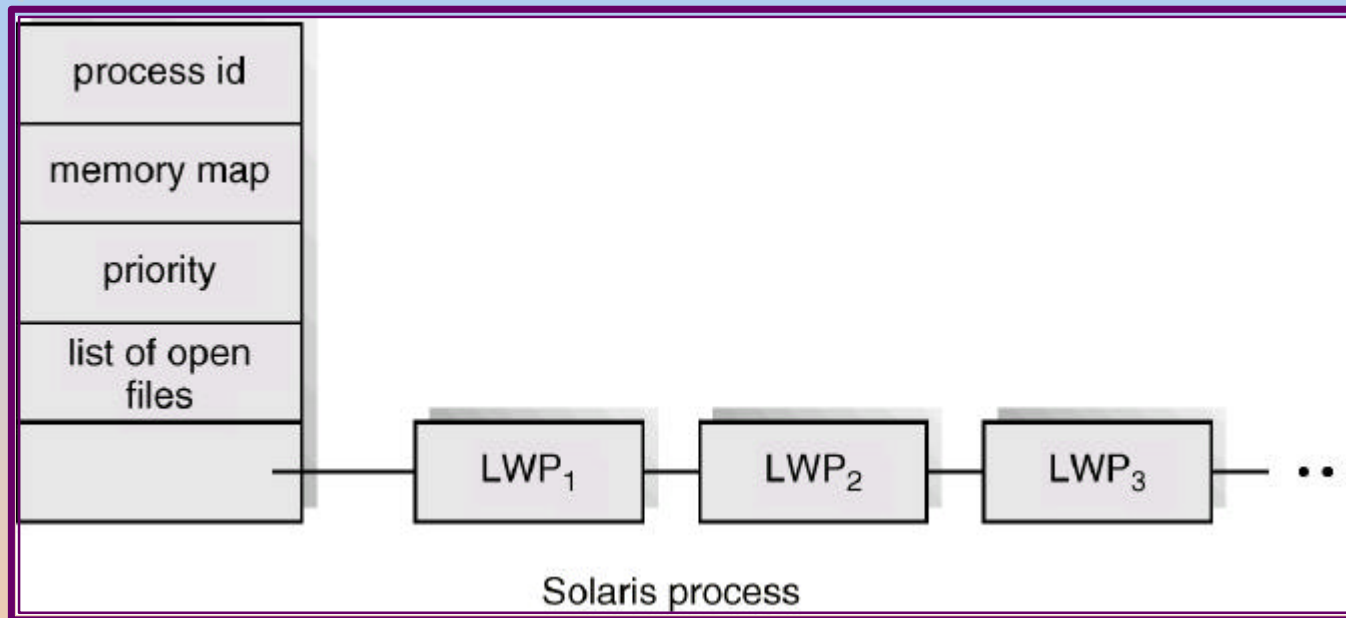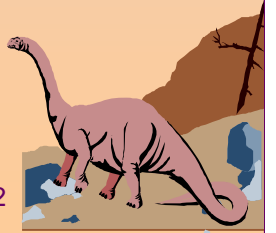
- Common in UNIX operating systems.

# Solaris 2 Threads



task 1     task 2     task 3

user-level thread

lightweight process

kernel thread

kernel

CPU

# Solaris Process



| | |
|---|---|
| process id | |
| memory map | |
| priority | |
| list of open files | |
| | LWP₁ — LWP₂ — LWP₃ — • • • |

Solaris process
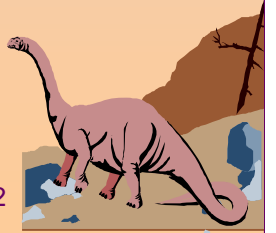
# Windows 2000 Threads

- Implements the one-to-one mapping.
- Each thread contains
  - a thread id
  - register set
  - separate user and kernel stacks
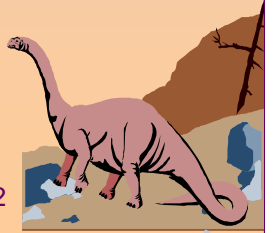  - private data storage area

# Linux Threads

- Linux refers to them as *tasks* rather than *threads*.
- Thread creation is done through clone() system call.
- Clone() allows a child task to share the address space of the parent task (process)
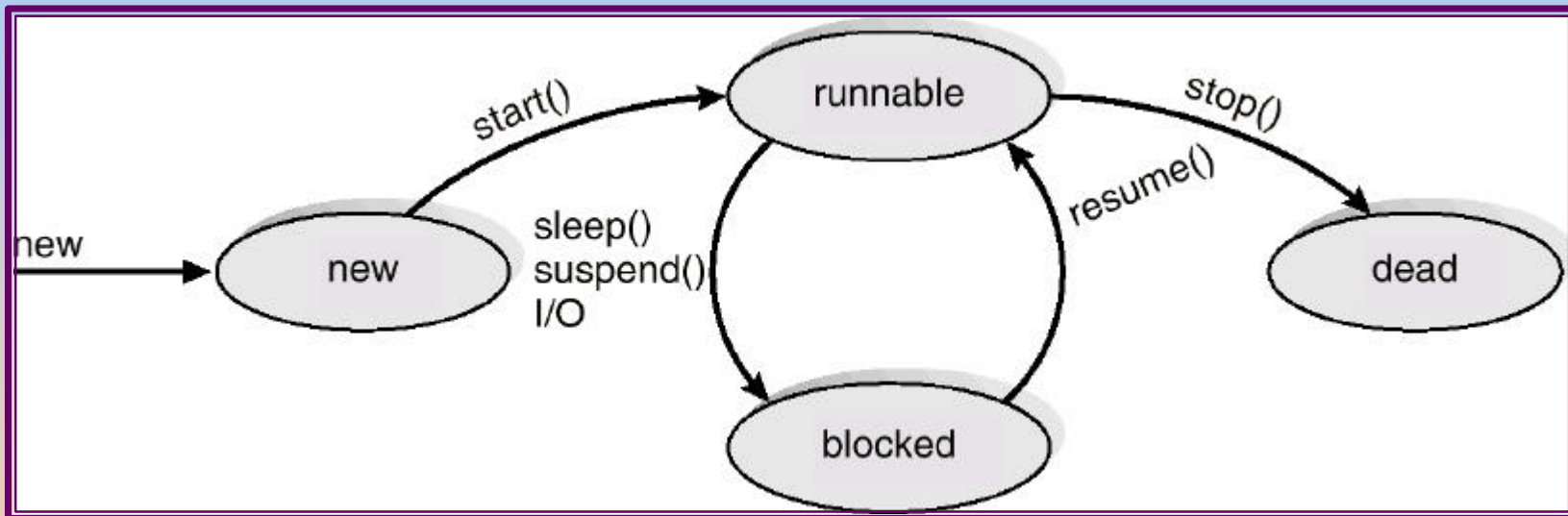
# Java Threads

■ Java threads may be created by:

    ✦ Extending Thread class

    ✦ Implementing the Runnable interface

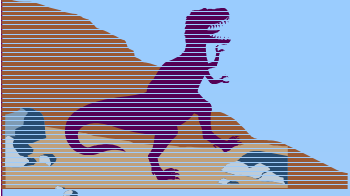■ Java threads are managed by the JVM.

# Java Thread States

# Chapter 6:  CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Multiple-Processor Scheduling
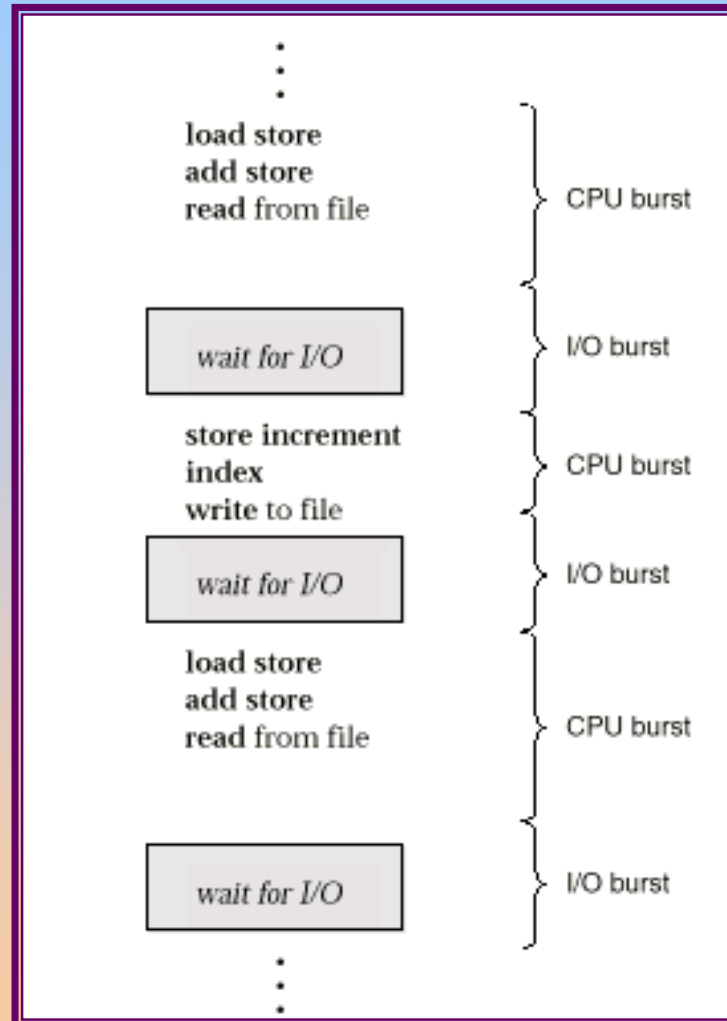- Real-Time Scheduling
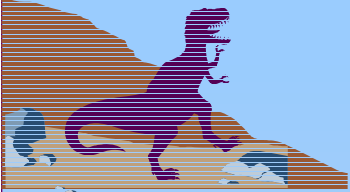- Algorithm Evaluation

# **Basic Concepts**

- Maximum CPU utilization obtained with multiprogramming

- CPU–I/O Burst Cycle – Process execution consists of a *cycle* of CPU execution and I/O wait.

- CPU burst distribution

# Alternating Sequence of CPU And I/O Bursts

```
        .
        .
        .
load store
add store
read from file      }  CPU burst


   wait for I/O      }  I/O burst


store increment
index
write to file        }  CPU burst


   wait for I/O      }  I/O burst


load store
add store
read from file      }  CPU burst


   wait for I/O      }  I/O burst
        .
        .
        .
```

# Histogram of CPU-burst Times

# CPU Scheduler

- Selects from among the processes in memory that are ready to execute, and allocates the CPU to one of them.

- CPU scheduling decisions may take place when a process:
    1. Switches from running to waiting state.
    2. Switches from running to ready state.
    3. Switches from waiting to ready.
    4. Terminates.

- Scheduling under 1 and 4 is *nonpreemptive*.
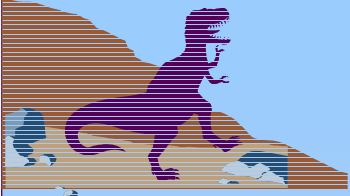
- All other scheduling is *preemptive.*

# Dispatcher

■ Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:

 ✦ switching context

 ✦ switching to user mode

 ✦ jumping to the proper location in the user program to restart that program

■ *Dispatch latency* – time it takes for the dispatcher to stop one process and start another running.
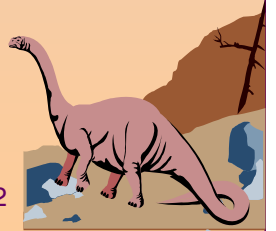
# Scheduling Criteria

- CPU utilization – keep the CPU as busy as possible

- Throughput – # of processes that complete their execution per time unit

- Turnaround time – amount of time to execute a particular process

- Waiting time – amount of time a process has been waiting in the ready queue

- Response time – amount of time it takes from when a request was submitted until the first response is produced, **not** output  (for time-sharing environment)
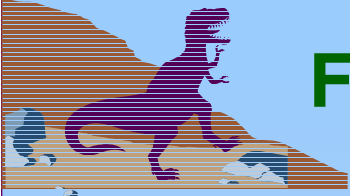
# **Optimization Criteria**

- Max CPU utilization
- Max throughput
- Min turnaround time
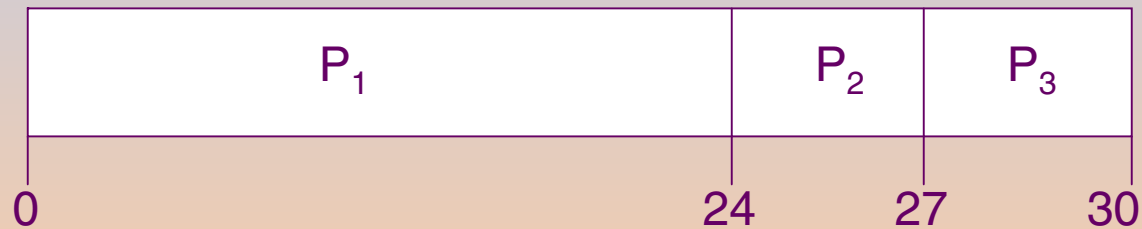- Min waiting time
- Min response time

# First-Come, First-Served (FCFS) Scheduling

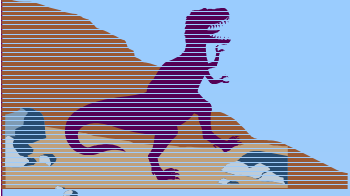| Process | Burst Time |
|---------|------------|
| $P_1$ | 24 |
| $P_2$ | 3 |
| $P_3$ | 3 |

- Suppose that the processes arrive in the order: $P_1$, $P_2$, $P_3$
  The Gantt Chart for the schedule is:

| P$_1$ | P$_2$ | P$_3$ |
|:---:|:---:|:---:|

0           24    27    30

- Waiting time for $P_1$ = 0; $P_2$ = 24; $P_3$ = 27
- Average waiting time: $(0 + 24 + 27)/3 = 17$

Suppose that the processes arrive in the order

$$P_2, P_3, P_1.$$

■ The Gantt chart for the schedule is:

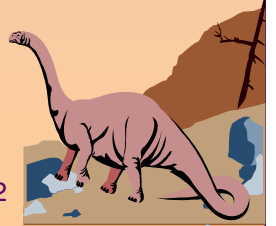| P$_2$ | P$_3$ | P$_1$ |
|---|---|---|

0    3    6                                    30

■ Waiting time for $P_1 = 6$; $P_2 = 0$; $P_3 = 3$

■ Average waiting time:   $(6 + 0 + 3)/3 = 3$

■ Much better than previous case.

■ *Convoy effect* short process behind long process

# Shortest-Job-First (SJR) Scheduling

- Associate with each process the length of its next CPU burst. Use these lengths to schedule the process with the shortest time.

- Two schemes:
  - nonpreemptive – once CPU given to the process it cannot be preempted until completes its CPU burst.
  - preemptive – if a new process arrives with CPU burst length less than remaining time of current executing process, preempt. This scheme is know as the Shortest-Remaining-Time-First (SRTF).

- SJF is optimal – gives minimum average waiting time for a given set of processes.

# Example of Non-Preemptive SJF

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| $P_1$ | 0.0 | 7 |
| $P_2$ | 2.0 | 4 |
| $P_3$ | 4.0 | 1 |
| $P_4$ | 5.0 | 4 |

- SJF (non-preemptive)

| P$_1$ | P$_3$ | P$_2$ | P$_4$ |
|-------|-------|-------|-------|

0    3       7  8       12      16

- Average waiting time = (0 + 6 + 3 + 7)/4 - 4

# Example of Preemptive SJF

| Process | Arrival Time | Burst Time |
|---------|-------------|-----------|
| $P_1$ | 0.0 | 7 |
| $P_2$ | 2.0 | 4 |
| $P_3$ | 4.0 | 1 |
| $P_4$ | 5.0 | 4 |

- SJF (preemptive)

| $P_1$ | $P_2$ | $P_3$ | $P_2$ | $P_4$ | $P_1$ |
|-------|-------|-------|-------|-------|-------|

0    2    4   5    7              11              16

- Average waiting time = (9 + 1 + 0 +2)/4 - 3

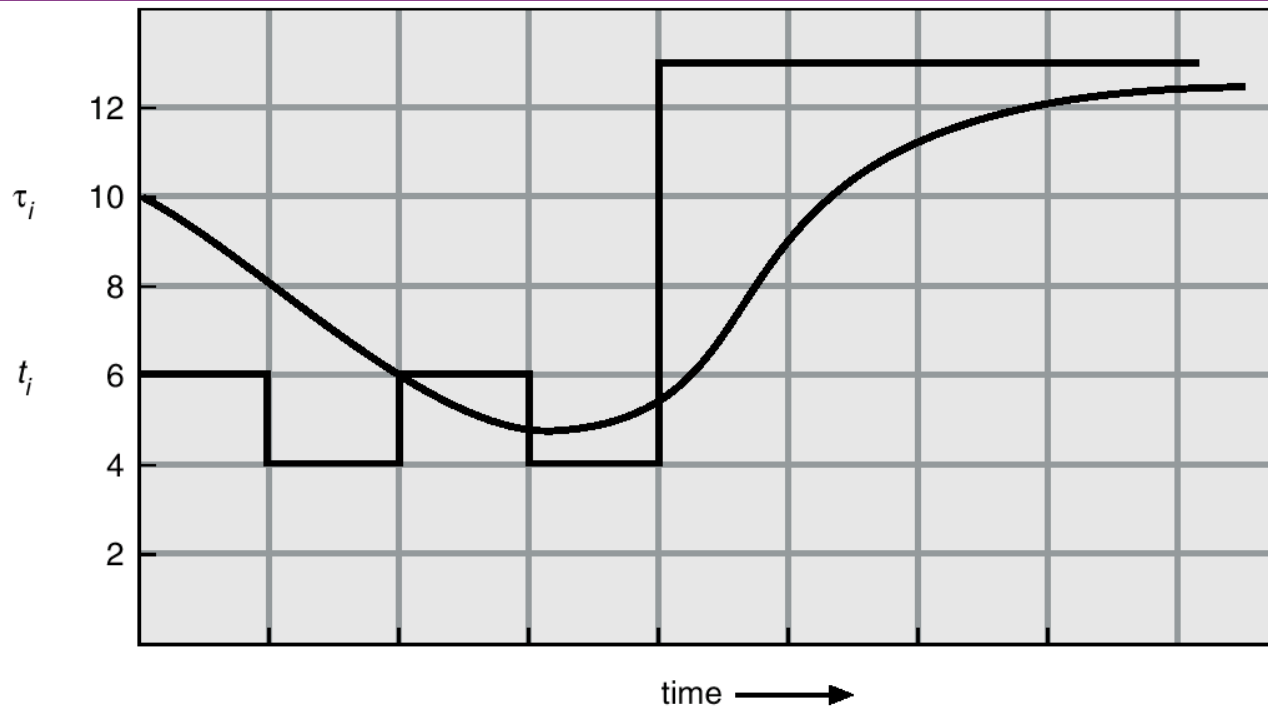# Determining Length of Next CPU Burst

■ Can only estimate the length.

■ Can be done by using the length of previous CPU bursts, using exponential averaging.

1. $t_n$ = actual lenght of $n^{th}$ CPU burst
2. $\tau_{n+1}$ = predicted value for the next CPU burst
3. $\alpha, 0 \le \alpha \le 1$
4. Define :

$$\tau_{n=1} = \alpha \, t_n + (1 - \alpha)\tau_n.$$

# Prediction of the Length of the Next CPU Burst



| CPU burst $(t_i)$ | | 6 | 4 | 6 | 4 | 13 | 13 | 13 | . . . |
|---|---|---|---|---|---|---|---|---|---|
| "guess" $(\tau_i)$ | 10 | 8 | 6 | 6 | 5 | 9 | 11 | 12 | . . . |

# Examples of Exponential Averaging

- $\alpha = 0$
  - $\tau_{n+1} = \tau_n$
  - Recent history does not count.
- $\alpha = 1$
  - $\tau_{n+1} = t_n$
  - Only the actual last CPU burst counts.
- If we expand the formula, we get:

$$\tau_{n+1} = \alpha\, t_n + (1 - \alpha)\, \alpha\, t_n - 1 + \dots$$
$$+ (1 - \alpha)^j\, \alpha\, t_n - 1 + \dots$$
$$+ (1 - \alpha)^{n=1}\, t_n\, \tau_0$$

- Since both $\alpha$ and $(1 - \alpha)$ are less than or equal to 1, each successive term has less weight than its predecessor.
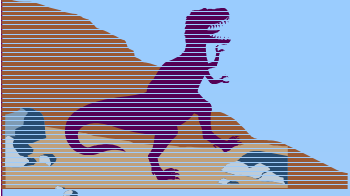
# Priority Scheduling

- A priority number (integer) is associated with each process

- The CPU is allocated to the process with the highest priority (smallest integer ≡ highest priority).
    - Preemptive
    - nonpreemptive

- SJF is a priority scheduling where priority is the predicted next CPU burst time.

- Problem ≡ Starvation – low priority processes may never execute.

- Solution ≡ Aging – as time progresses increase the priority of the process.

# Round Robin (RR)

- Each process gets a small unit of CPU time (*time quantum*), usually 10-100 milliseconds.  After this time has elapsed, the process is preempted and added to the end of the ready queue.

- If there are *n* processes in the ready queue and the time quantum is *q*, then each process gets 1/*n* of the CPU time in chunks of at most *q* time units at once.  No process waits more than (*n*-1)*q* time units.

- Performance
  - *q* large $\Rightarrow$ FIFO
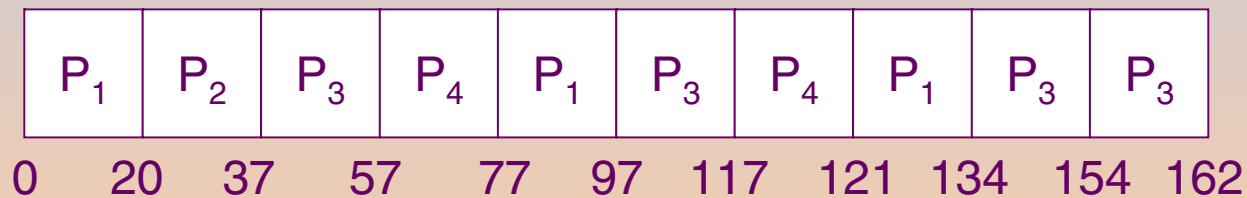  - *q* small $\Rightarrow$ *q* must be large with respect to context switch, otherwise overhead is too high.

# Example of RR with Time Quantum = 20

| Process | Burst Time |
|---------|------------|
| $P_1$ | 53 |
| $P_2$ | 17 |
| $P_3$ | 68 |
| $P_4$ | 24 |

■ The Gantt chart is:

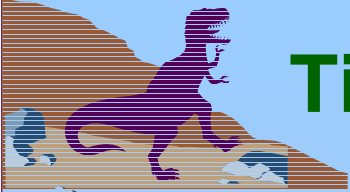| $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_1$ | $P_3$ | $P_4$ | $P_1$ | $P_3$ | $P_3$ |
|---|---|---|---|---|---|---|---|---|---|

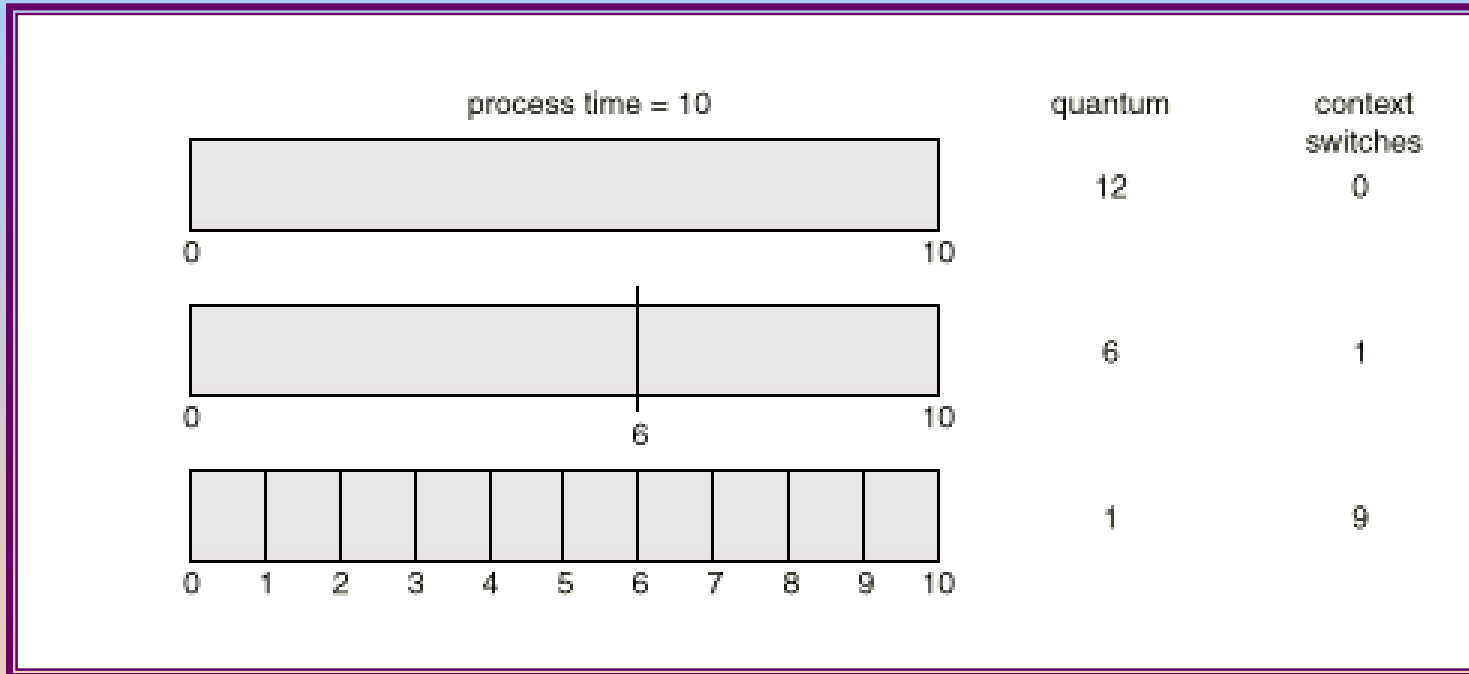0    20    37    57    77    97    117    121    134    154    162

■ Typically, higher average turnaround than SJF, but better *response*.

# Time Quantum and Context Switch Time

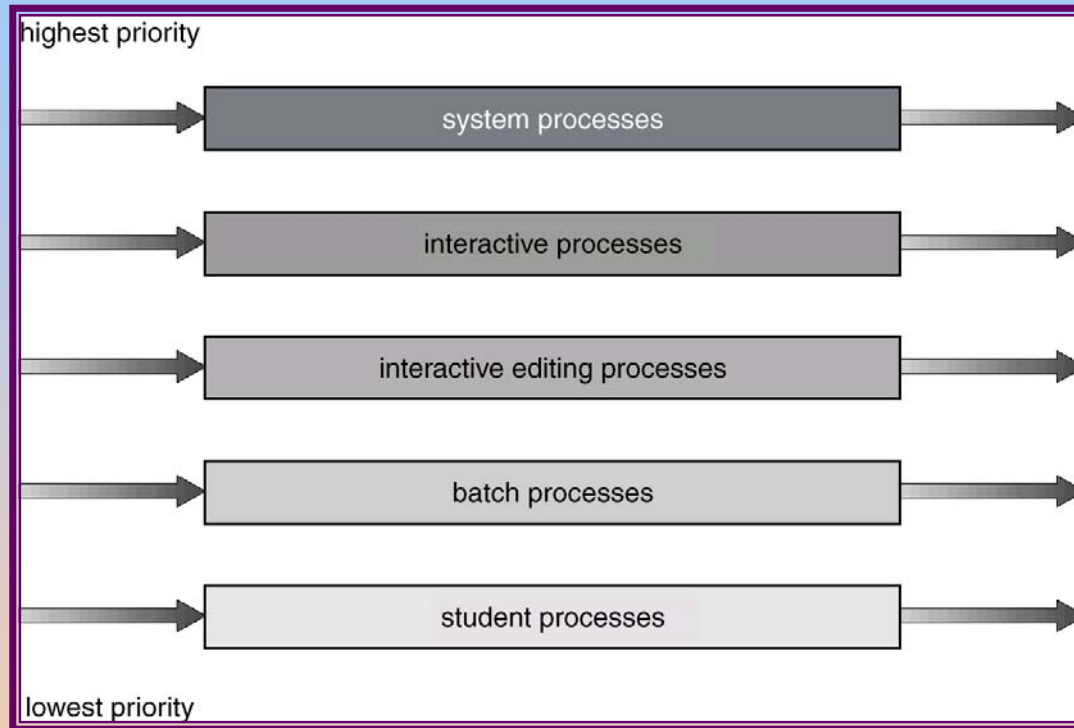# Turnaround Time Varies With The Time Quantum



| process | time |
|---------|------|
| $P_1$   | 6    |
| $P_2$   | 3    |
| $P_3$   | 1    |
| $P_4$   | 7    |

# Multilevel Queue

- Ready queue is partitioned into separate queues:
  foreground (interactive)
  background (batch)

- Each queue has its own scheduling algorithm,
  foreground – RR
  background – FCFS

- Scheduling must be done between the queues.

  - Fixed priority scheduling; (i.e., serve all from foreground then from background). Possibility of starvation.

  - Time slice – each queue gets a certain amount of CPU time which it can schedule amongst its processes; i.e., 80% to foreground in RR

  - 20% to background in FCFS

# Multilevel Queue Scheduling



highest priority

→ system processes →

→ interactive processes →

→ interactive editing processes →

→ batch processes →

→ student processes →

lowest priority

# Multilevel Feedback Queue

- A process can move between the various queues; aging can be implemented this way.

- Multilevel-feedback-queue scheduler defined by the following parameters:
  - number of queues
  - scheduling algorithms for each queue
  - method used to determine when to upgrade a process
  - method used to determine when to demote a process
  - method used to determine which queue a process will enter when that process needs service
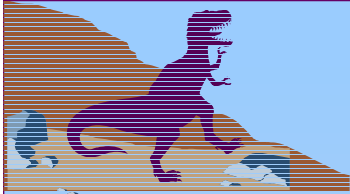
# Example of Multilevel Feedback Queue

- Three queues:
  - ✦ $Q_0$ – time quantum 8 milliseconds
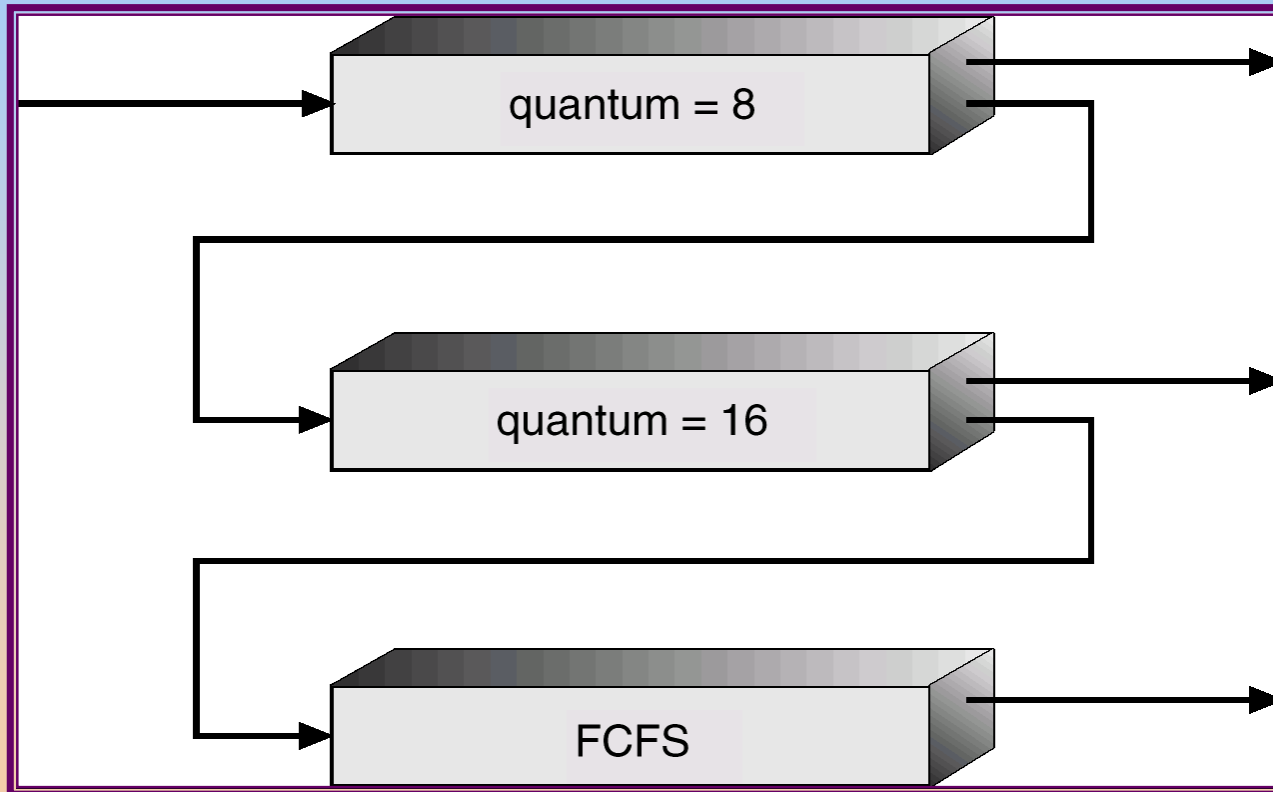  - ✦ $Q_1$ – time quantum 16 milliseconds
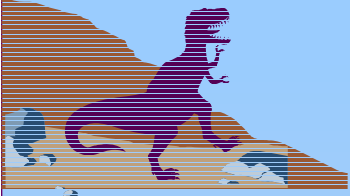  - ✦ $Q_2$ – FCFS
- Scheduling
  - ✦ A new job enters queue $Q_0$ which is served FCFS. When it gains CPU, job receives 8 milliseconds.  If it does not finish in 8 milliseconds, job is moved to queue $Q_1$.
  - ✦ At $Q_1$ job is again served FCFS and receives 16 additional milliseconds.  If it still does not complete, it is preempted and moved to queue $Q_2$.

# Multilevel Feedback Queues

# Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available.

- *Homogeneous processors* within a multiprocessor.

- *Load sharing*

- *Asymmetric multiprocessing* – only one processor accesses the system data structures, alleviating the need for data sharing.
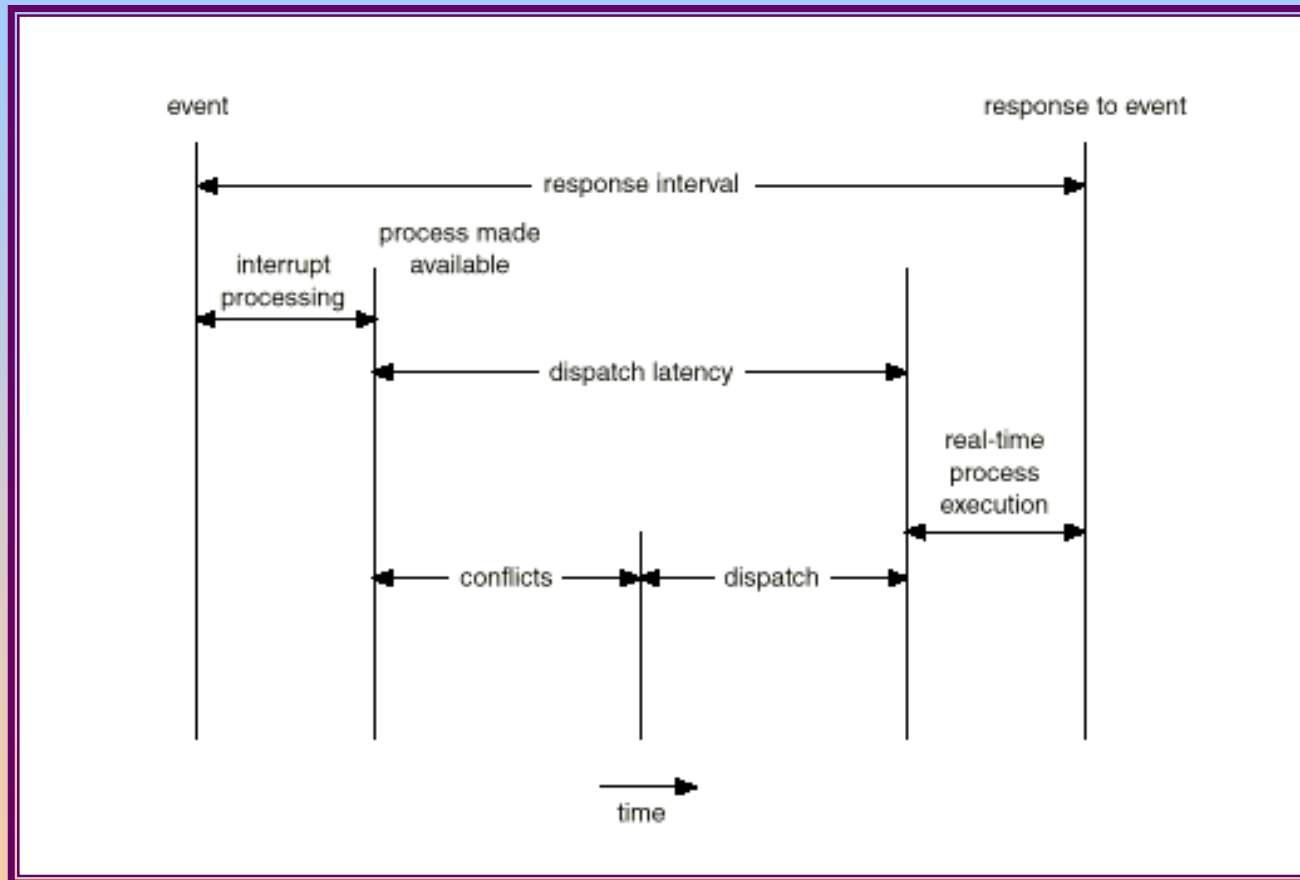
# Real-Time Scheduling

■ *Hard real-time* systems – required to complete a critical task within a guaranteed amount of time.

■ *Soft real-time* computing – requires that critical processes receive priority over less fortunate ones.
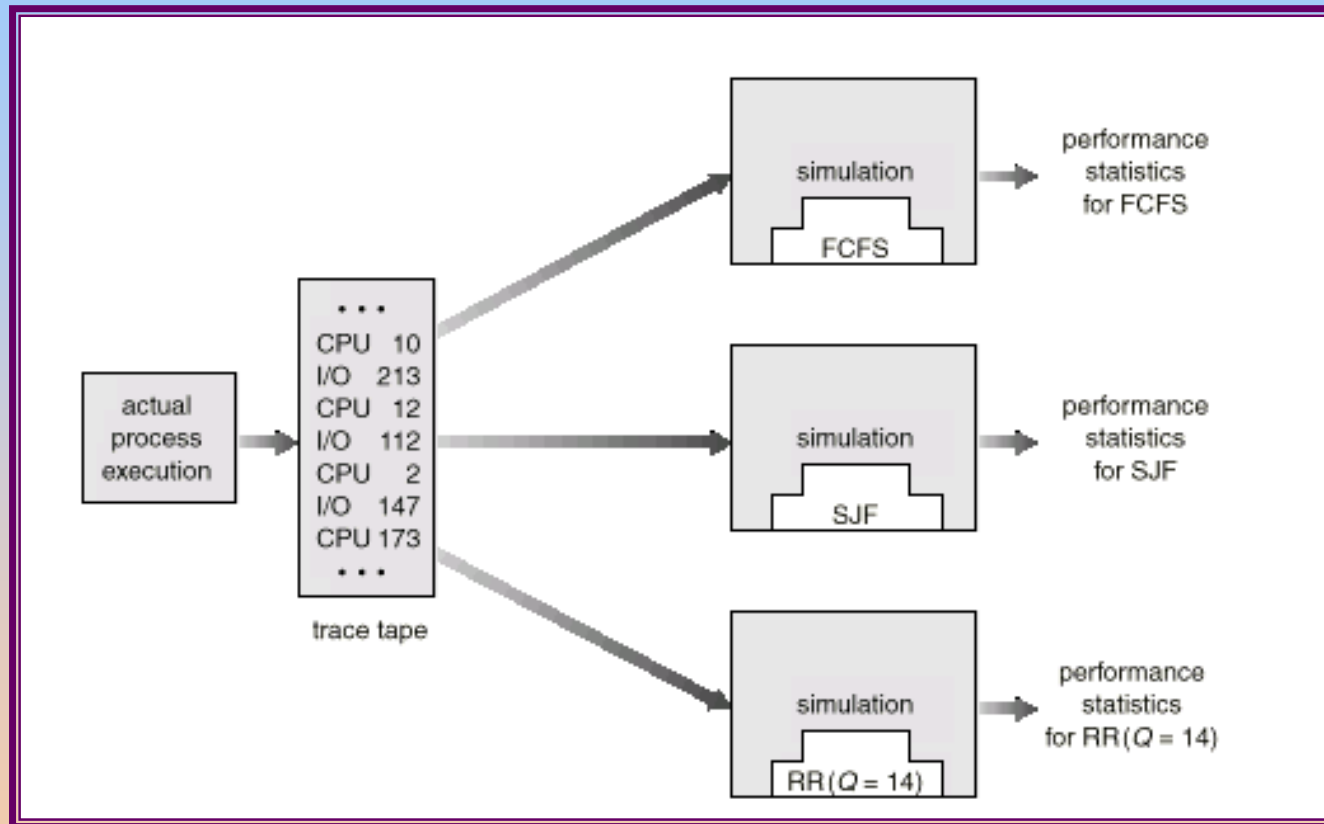
# Dispatch Latency

# Algorithm Evaluation

- Deterministic modeling – takes a particular predetermined workload and defines the performance of each algorithm for that workload.
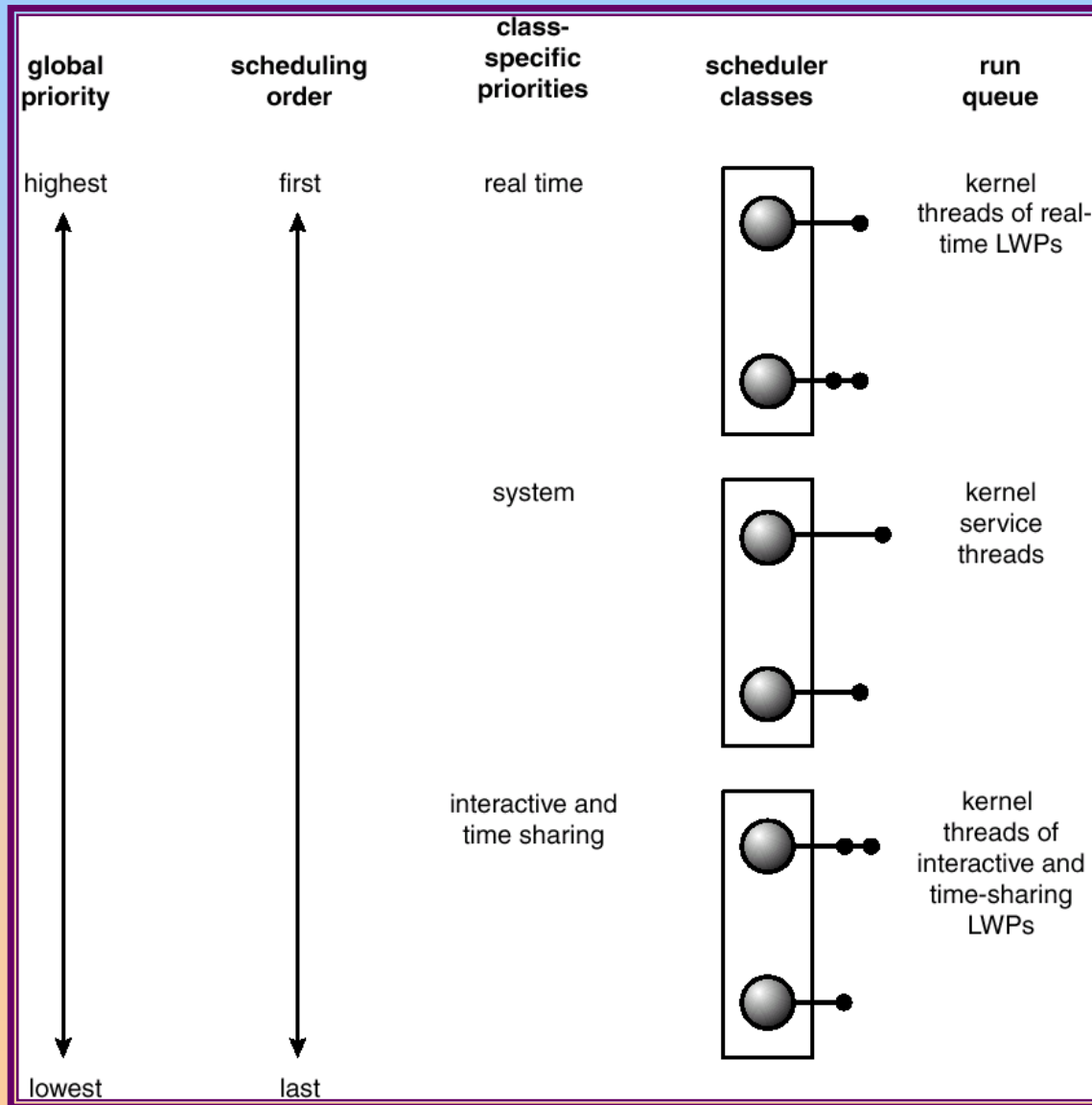
- Queueing models

- Implementation

# Evaluation of CPU Schedulers by Simulation

# Solaris 2 Scheduling

# Windows 2000 Priorities

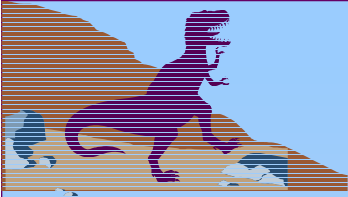|  | real-time | high | above normal | normal | below normal | idle priority |
|---|---|---|---|---|---|---|
| time-critical | 31 | 15 | 15 | 15 | 15 | 15 |
| highest | 26 | 15 | 12 | 10 | 8 | 6 |
| above normal | 25 | 14 | 11 | 9 | 7 | 5 |
| normal | 24 | 13 | 10 | 8 | 6 | 4 |
| below normal | 23 | 12 | 9 | 7 | 5 | 3 |
| lowest | 22 | 11 | 8 | 6 | 4 | 2 |
| idle | 16 | 1 | 1 | 1 | 1 | 1 |

# Chapter 7:  Process Synchronization

- Background
- The Critical-Section Problem
- Synchronization Hardware
- Semaphores
- Classical Problems of Synchronization
- Critical Regions
- Monitors
- Synchronization in Solaris 2 & Windows 2000

# Background

- Concurrent access to shared data may result in data inconsistency.

- Maintaining data consistency requires mechanisms to ensure the orderly execution of cooperating processes.

- Shared-memory solution to bounded-butter problem (Chapter 4) allows at most $n-1$ items in buffer at the same time. A solution, where all $N$ buffers are used is not simple.

  - Suppose that we modify the producer-consumer code by adding a variable *counter*, initialized to 0 and incremented each time a new item is added to the buffer

# Bounded-Buffer

- Shared data

```
#define BUFFER_SIZE 10
typedef struct {
   . . .
} item;
item buffer[BUFFER_SIZE];
int in = 0;
int out = 0;
int counter = 0;
```
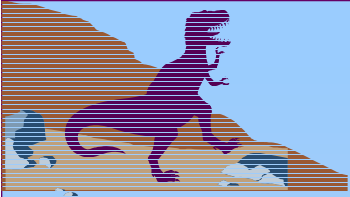
# Bounded-Buffer

■ Producer process

```
item nextProduced;

while (1) {
        while (counter == BUFFER_SIZE)
                ; /* do nothing */
        buffer[in] = nextProduced;
        in = (in + 1) % BUFFER_SIZE;
        counter++;
}
```

# Bounded-Buffer

- Consumer process

```
item nextConsumed;

while (1) {
        while (counter == 0)
                ; /* do nothing */
        nextConsumed = buffer[out];
        out = (out + 1) % BUFFER_SIZE;
        counter--;
}
```

# Bounded Buffer

■ The statements

**counter++;**
**counter--;**

must be performed *atomically*.

■ Atomic operation means an operation that completes in its entirety without interruption.

# Bounded Buffer

■ The statement "**count++**" may be implemented in machine language as:

**register1 = counter**
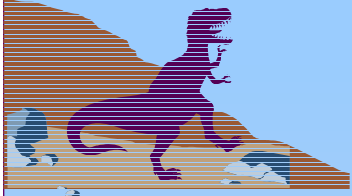
**register1 = register1 + 1**
**counter = register1**

■ The statement "**count—**" may be implemented as:

**register2 = counter**
**register2 = register2 – 1**
**counter = register2**

# **Bounded Buffer**

- If both the producer and consumer attempt to update the buffer concurrently, the assembly language statements may get interleaved.

- Interleaving depends upon how the producer and consumer processes are scheduled.

# Bounded Buffer

■ Assume **counter** is initially 5. One interleaving of statements is:

producer: **register1 = counter** (*register1 = 5*)
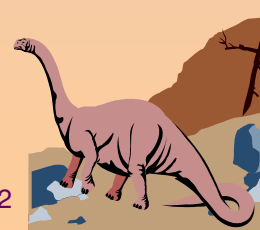producer: **register1 = register1 + 1** (*register1 = 6*)
consumer: **register2 = counter** (*register2 = 5*)
consumer: **register2 = register2 – 1** (*register2 = 4*)
producer: **counter = register1** (*counter = 6*)
consumer: **counter = register2** (*counter = 4*)

■ The value of **count** may be either 4 or 6, where the correct result should be 5.

# Race Condition

- **Race condition**: The situation where several processes access – and manipulate shared data concurrently. The final value of the shared data depends upon which process finishes last.

- To prevent race conditions, concurrent processes must be **synchronized**.
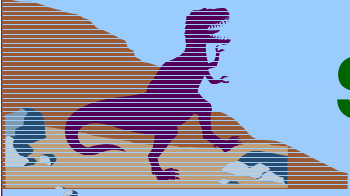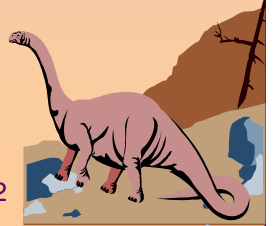
# The Critical-Section Problem

- *n* processes all competing to use some shared data

- Each process has a code segment, called *critical section*, in which the shared data is accessed.

- Problem – ensure that when one process is executing in its critical section, no other process is allowed to execute in its critical section.

# **Solution to Critical-Section Problem**

1. **Mutual Exclusion**. If process $P_i$ is executing in its critical section, then no other processes can be executing in their critical sections.

2. **Progress**. If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the processes that will enter the critical section next cannot be postponed indefinitely.

3. **Bounded Waiting**. A bound must exist on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted.

   - Assume that each process executes at a nonzero speed
   - No assumption concerning relative speed of the $n$ processes.

# Initial Attempts to Solve Problem

- Only 2 processes, $P_0$ and $P_1$
- General structure of process $P_i$ (other process $P_j$)

> **do** {
>
>      entry section
>
>         critical section
>
>      exit section
>
>         reminder section
>
> } **while (1)**;

- Processes may share some common variables to synchronize their actions.

# Algorithm 1

- Shared variables:
    - **int turn**;
      initially **turn = 0**
    - **turn - i** $\Rightarrow$ $P_i$ can enter its critical section
- Process $P_i$
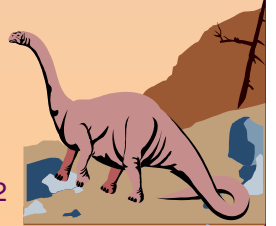
        **do** {
                **while (turn != i)** ;
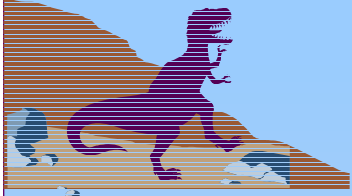                    critical section
                **turn = j**;
                    reminder section
        } **while (1)**;

- Satisfies mutual exclusion, but not progress

# Algorithm 2

- Shared variables
  - **boolean flag[2]**;
    initially **flag [0] = flag [1] = false.**
  - **flag [i] = true** $\Rightarrow P_i$ ready to enter its critical section
- Process $P_i$

```
do {
    flag[i] := true;
    while (flag[j]) ;
        critical section
    flag [i] = false;
        remainder section
} while (1);
```

- Satisfies mutual exclusion, but not progress requirement.
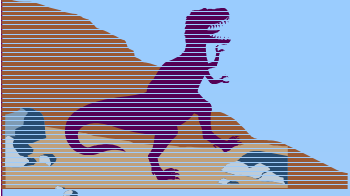
# Algorithm 3

- Combined shared variables of algorithms 1 and 2.
- Process $P_i$

```
do {
        flag [i]:= true;
        turn = j;
        while (flag [j] and turn = j) ;
            critical section
        flag [i] = false;
            remainder section
} while (1);
```

- Meets all three requirements; solves the critical-section problem for two processes.

# Bakery Algorithm

Critical section for n processes

- Before entering its critical section, process receives a number. Holder of the smallest number enters the critical section.

- If processes $P_i$ and $P_j$ receive the same number, if $i < j$, then $P_i$ is served first; else $P_j$ is served first.

- The numbering scheme always generates numbers in increasing order of enumeration; i.e., 1,2,3,3,3,3,4,5...

# Bakery Algorithm

- Notation $\leq \equiv$ lexicographical order (ticket #, process id #)
  - ◆ $(a,b) < c,d)$ if $a < c$ or if $a = c$ and $b < d$
  - ◆ max $(a_0,\ldots, a_{n-1})$ is a number, $k$, such that $k \geq a_i$ for $i$ - 0, $\ldots$, $n - 1$
- Shared data

  **boolean choosing[n];**

  **int number[n];**

  Data structures are initialized to **false** and **0** respectively

# Bakery Algorithm

```
do {
    choosing[i] = true;
    number[i] = max(number[0], number[1], …, number [n – 1])+1;
    choosing[i] = false;
    for (j = 0; j < n; j++) {
            while (choosing[j]) ;
            while ((number[j] != 0) && (number[j,j] < number[i,i])) ;
    }
        critical section
    number[i] = 0;
        remainder section
} while (1);
```

# Synchronization Hardware

- Test and modify the content of a word atomically

```
boolean TestAndSet(boolean &target) {
    boolean rv = target;
    tqrget = true;

    return rv;
}
```
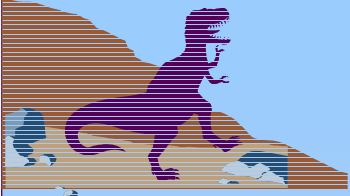
# Mutual Exclusion with Test-and-Set

- Shared data:

  **boolean lock = false;**

- Process $P_i$

  **do {**

      **while (TestAndSet(lock)) ;**

        critical section

      **lock = false;**

        remainder section

  **}**

# Synchronization Hardware

■ Atomically swap two variables.

```
void Swap(boolean &a, boolean &b) {
    boolean temp = a;
    a = b;
    b = temp;
}
```

# Mutual Exclusion with Swap

- Shared data (initialized to **false**):
  **boolean lock;**
  **boolean waiting[n];**

- Process $P_i$
  **do {**
      **key = true;**
      **while (key == true)**
           **Swap(lock,key);**
      critical section
      **lock = false;**
      remainder section
  **}**

# Semaphores

- Synchronization tool that does not require busy waiting.
- Semaphore $S$ – integer variable
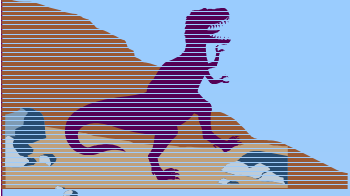- can only be accessed via two indivisible (atomic) operations

> *wait* ($S$):
>
> **while $S \leq 0$ do *no-op*;**
> **$S$--;**
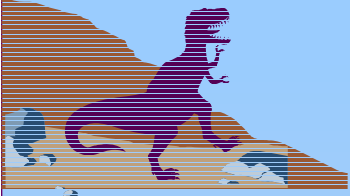>
>
> *signal* ($S$):
>
> **$S$++;**

# Critical Section of *n* Processes

- Shared data:

  **semaphore mutex;** *//*initially *mutex* = 1

- Process *Pi:*

  ```
  do {
      wait(mutex);
          critical section
      signal(mutex);
          remainder section
  } while (1);
  ```

# Semaphore Implementation

- Define a semaphore as a record

  **typedef struct {**

  **int value;**
  **struct process *L;**
  **} semaphore;**

- Assume two simple operations:
  - **block** suspends the process that invokes it.
  - **wakeup(P)** resumes the execution of a blocked process **P**.

# Implementation

- Semaphore operations now defined as

  *wait*(*S*):

  **S.value--;**
  **if (S.value < 0) {**

      add this process to **S.L;**
      **block;**

  **}**

  *signal*(*S*):

  **S.value++;**
  **if (S.value <= 0) {**

      remove a process **P** from **S.L;**
      **wakeup(P);**

  **}**

# Semaphore as a General Synchronization Tool

- Execute $B$ in $P_j$ only after $A$ executed in $P_i$
- Use semaphore *flag* initialized to 0
- Code:

$$P_i$$

$\vdots$

$A$

*signal*(*flag*)

$$P_j$$

$\vdots$

*wait*(*flag*)

$B$

# Deadlock and Starvation

- **Deadlock** – two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes.

- Let $S$ and $Q$ be two semaphores initialized to 1

| $P_0$ | $P_1$ |
|---|---|
| wait($S$); | wait($Q$); |
| wait($Q$); | wait($S$); |
| $\vdots$ | $\vdots$ |
| signal($S$); | signal($Q$); |
| signal($Q$) | signal($S$); |

- **Starvation** – indefinite blocking.  A process may never be removed from the semaphore queue in which it is suspended.

# Two Types of Semaphores

- *Counting* semaphore – integer value can range over an unrestricted domain.

- *Binary* semaphore – integer value can range only between 0 and 1; can be simpler to implement.

- Can implement a counting semaphore $S$ as a binary semaphore.

# Implementing *S* as a Binary Semaphore

- Data structures:

  **binary-semaphore S1, S2;**

  **int C:**

- Initialization:

  **S1 = 1**

  **S2 = 0**

  **C =** initial value of semaphore **S**

# Implementing $S$

- *wait* operation

```
wait(S1);
C--;
if (C < 0) {
            signal(S1);
            wait(S2);
}
signal(S1);
```

- *signal* operation

```
wait(S1);
C ++;
if (C <= 0)
        signal(S2);
else
        signal(S1);
```

# Classical Problems of Synchronization

- Bounded-Buffer Problem

- Readers and Writers Problem

- Dining-Philosophers Problem

# Bounded-Buffer Problem

- Shared data

  **semaphore full, empty, mutex;**

  Initially:

  **full = 0, empty = n, mutex = 1**

# Bounded-Buffer Problem Producer Process

```
do {
        …
    produce an item in nextp
        …
    wait(empty);
    wait(mutex);
        …
    add nextp to buffer
        …
    signal(mutex);
    signal(full);
} while (1);
```

# Bounded-Buffer Problem Consumer Process

```
do {
    wait(full)
    wait(mutex);

        …
    remove an item from buffer to nextc

        …
    signal(mutex);
    signal(empty);

        …
    consume the item in nextc

        …
} while (1);
```

# Readers-Writers Problem

■ Shared data

**semaphore mutex, wrt;**

Initially

**mutex = 1, wrt = 1, readcount = 0**

# Readers-Writers Problem Writer Process

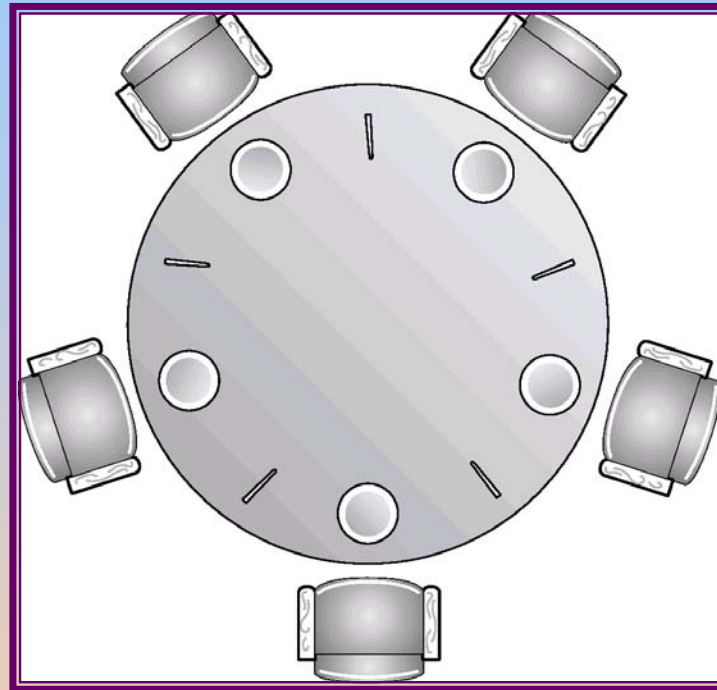**wait(wrt);**

**…**

writing is performed

**…**

**signal(wrt);**

# Readers-Writers Problem Reader Process

```
wait(mutex);
readcount++;
if (readcount == 1)
        wait(rt);
signal(mutex);

        …
    reading is performed
        …
wait(mutex);
readcount--;
if (readcount == 0)
    signal(wrt);
signal(mutex):
```

# Dining-Philosophers Problem



- Shared data

    **semaphore chopstick[5];**

Initially all values are 1

# Dining-Philosophers Problem

- Philosopher *i*:

```
do {
    wait(chopstick[i])
    wait(chopstick[(i+1) % 5])

      …
     eat

      …
    signal(chopstick[i]);
    signal(chopstick[(i+1) % 5]);

      …
     think

      …
    } while (1);
```

# Critical Regions

- High-level synchronization construct

- A shared variable **v** of type **T**, is declared as:

    **v: shared T**

- Variable **v** accessed only inside statement

    **region v when B do S**

    where **B** is a boolean expression.

- While statement **S** is being executed, no other process can access variable **v**.

# Critical Regions

- Regions referring to the same shared variable exclude each other in time.

- When a process tries to execute the region statement, the Boolean expression $B$ is evaluated. If $B$ is true, statement $S$ is executed. If it is false, the process is delayed until $B$ becomes true and no other process is in the region associated with $v$.

# Example – Bounded Buffer

- Shared data:

  **struct buffer {**
        **int pool[n];**
        **int count, in, out;**
  **}**

# Bounded Buffer Producer Process

- Producer process inserts **nextp** into the shared buffer

```
region buffer when( count < n) {
        pool[in] = nextp;
        in:= (in+1) % n;
        count++;
}
```

# Bounded Buffer Consumer Process

■ Consumer process removes an item from the shared buffer and puts it in **nextc**

```
region buffer when (count > 0) {
    nextc = pool[out];
    out = (out+1) % n;
    count--;
}
```

# Implementation region *x* when *B* do *S*

- Associate with the shared variable *x*, the following variables:

  **semaphore mutex, first-delay, second-delay;**
  **int first-count, second-count;**

- Mutually exclusive access to the critical section is provided by **mutex**.

- If a process cannot enter the critical section because the Boolean expression **B** is false, it initially waits on the **first-delay** semaphore; moved to the **second-delay** semaphore before it is allowed to reevaluate *B*.

# Implementation

- Keep track of the number of processes waiting on **first-delay** and **second-delay**, with **first-count** and **second-count** respectively.

- The algorithm assumes a FIFO ordering in the queuing of processes for a semaphore.

- For an arbitrary queuing discipline, a more complicated implementation is required.

# Monitors

■ High-level synchronization construct that allows the safe sharing of an abstract data type among concurrent processes.

```
monitor monitor-name
{
        shared variable declarations
        procedure body P1 (...) {
           . . .
        }
        procedure body P2 (...) {
           . . .
        }
        procedure body Pn (...) {
            . . .
        }
        {
           initialization code
        }
}
```

# Monitors

- To allow a process to wait within the monitor, a **condition** variable must be declared, as

  **condition x, y;**

- Condition variable can only be used with the operations **wait** and **signal**.

  - The operation

    **x.wait();**

    means that the process invoking this operation is suspended until another process invokes

    **x.signal();**

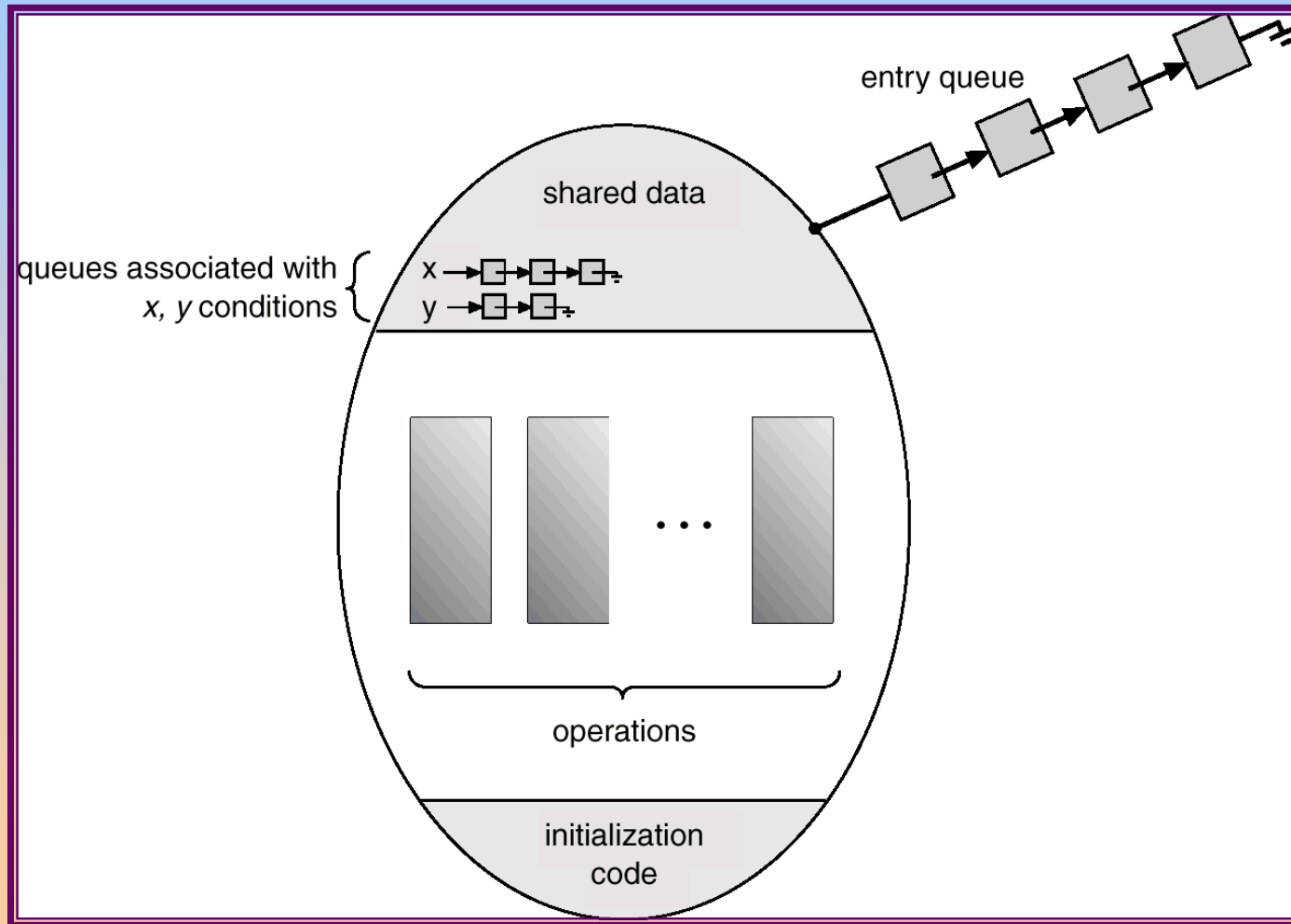  - The **x.signal** operation resumes exactly one suspended process. If no process is suspended, then the **signal** operation has no effect.

# Schematic View of a Monitor



entry queue

shared data

... 

operations

initialization code

# Monitor With Condition Variables

# Dining Philosophers Example

```
monitor dp
{
  enum {thinking, hungry, eating} state[5];
  condition self[5];
  void pickup(int i)          // following slides
  void putdown(int i)         // following slides
  void test(int i)            // following slides
  void init() {
      for (int i = 0; i < 5; i++)
              state[i] = thinking;
  }
}
```

# Dining Philosophers

```
void pickup(int i) {
        state[i] = hungry;
        test[i];
        if (state[i] != eating)
                self[i].wait();
}

void putdown(int i) {
        state[i] = thinking;
        // test left and right neighbors
        test((i+4) % 5);
        test((i+1) % 5);
}
```

# Dining Philosophers

```
void test(int i) {
        if ( (state[(I + 4) % 5] != eating) &&
          (state[i] == hungry) &&
          (state[(i + 1) % 5] != eating)) {
                state[i] = eating;
                self[i].signal();
        }
}
```

# Monitor Implementation Using Semaphores

- Variables

  **semaphore mutex;  // (initially  = 1)**
  **semaphore next;     // (initially  = 0)**
  **int next-count = 0;**

- Each external procedure *F* will be replaced by

  **wait(mutex);**

  …

  body of *F*;

  …

  **if (next-count > 0)**
  **   signal(next)**
  **else**
  **   signal(mutex);**

- Mutual exclusion within a monitor is ensured.

# Monitor Implementation

■ For each condition variable **x**, we have:

> **semaphore x-sem; // (initially = 0)**
> **int x-count = 0;**

■ The operation **x.wait** can be implemented as:

> **x-count++;**
> **if (next-count > 0)**
> **signal(next);**
> **else**
> **signal(mutex);**
> **wait(x-sem);**
> **x-count--;**

# Monitor Implementation

- The operation **x.signal** can be implemented as:

```
if (x-count > 0) {
        next-count++;
        signal(x-sem);
        wait(next);
        next-count--;
}
```

# Monitor Implementation

- *Conditional-wait* construct: **x.wait(c);**
  - **c** – integer expression evaluated when the **wait** operation is executed.
  - value of **c** (a *priority number*) stored with the name of the process that is suspended.
  - when **x.signal** is executed, process with smallest associated priority number is resumed next.
- Check two conditions to establish correctness of system:
  - User processes must always make their calls on the monitor in a correct sequence.
  - Must ensure that an uncooperative process does not ignore the mutual-exclusion gateway provided by the monitor, and try to access the shared resource directly, without using the access protocols.

# Solaris 2 Synchronization

- Implements a variety of locks to support multitasking, multithreading (including real-time threads), and multiprocessing.

- Uses *adaptive mutexes* for efficiency when protecting data from short code segments.

- Uses *condition variables* and *readers-writers* locks when longer sections of code need access to data.

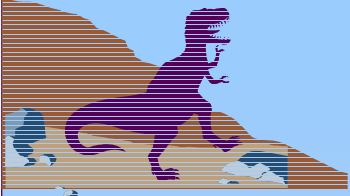- Uses *turnstiles* to order the list of threads waiting to acquire either an adaptive mutex or reader-writer lock.

# Windows 2000 Synchronization

- Uses interrupt masks to protect access to global resources on uniprocessor systems.

- Uses *spinlocks* on multiprocessor systems.

- Also provides *dispatcher objects* which may act as wither mutexes and semaphores.

- Dispatcher objects may also provide *events*. An event acts much like a condition variable.

# Chapter 8: Deadlocks

- System Model
- Deadlock Characterization
- Methods for Handling Deadlocks
- Deadlock Prevention
- Deadlock Avoidance
- Deadlock Detection
- Recovery from Deadlock
- Combined Approach to Deadlock Handling

# The Deadlock Problem

- A set of blocked processes each holding a resource and waiting to acquire a resource held by another process in the set.

- Example
  - System has 2 tape drives.
  - $P_1$ and $P_2$ each hold one tape drive and each needs another one.

- Example
  - semaphores $A$ and $B$, initialized to 1

|  $P_0$ | $P_1$ |
|---|---|
| *wait (A);* | *wait(B)* |
| *wait (B);* | *wait(A)* |

# Bridge Crossing Example

- Traffic only in one direction.

- Each section of a bridge can be viewed as a resource.

- If a deadlock occurs, it can be resolved if one car backs up (preempt resources and rollback).

- Several cars may have to be backed up if a deadlock occurs.

- Starvation is possible.

# System Model

- Resource types $R_1, R_2, \ldots, R_m$
    - *CPU cycles, memory space, I/O devices*
- Each resource type $R_i$ has $W_i$ instances.
- Each process utilizes a resource as follows:
    - ✦ request
    - ✦ use
    - ✦ release

# Deadlock Characterization

Deadlock can arise if four conditions hold simultaneously.

- **Mutual exclusion:** only one process at a time can use a resource.

- **Hold and wait:** a process holding at least one resource is waiting to acquire additional resources held by other processes.

- **No preemption:** a resource can be released only voluntarily by the process holding it, after that process has completed its task.

- **Circular wait:** there exists a set $\{P_0, P_1, \ldots, P_0\}$ of waiting processes such that $P_0$ is waiting for a resource that is held by $P_1$, $P_1$ is waiting for a resource that is held by $P_2, \ldots, P_{n-1}$ is waiting for a resource that is held by $P_n$, and $P_0$ is waiting for a resource that is held by $P_0$.

# Resource-Allocation Graph

A set of vertices $V$ and a set of edges $E$.

- V is partitioned into two types:
    - $P = \{P_1, P_2, \ldots, P_n\}$, the set consisting of all the processes in the system.

    - $R = \{R_1, R_2, \ldots, R_m\}$, the set consisting of all resource types in the system.
- request edge – directed edge $P_1 \rightarrow R_j$
- assignment edge – directed edge $R_j \rightarrow P_i$
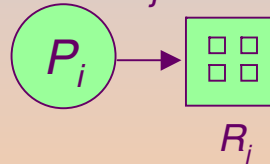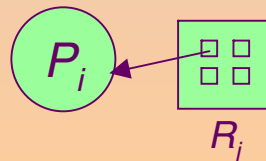
# Resource-Allocation Graph (Cont.)
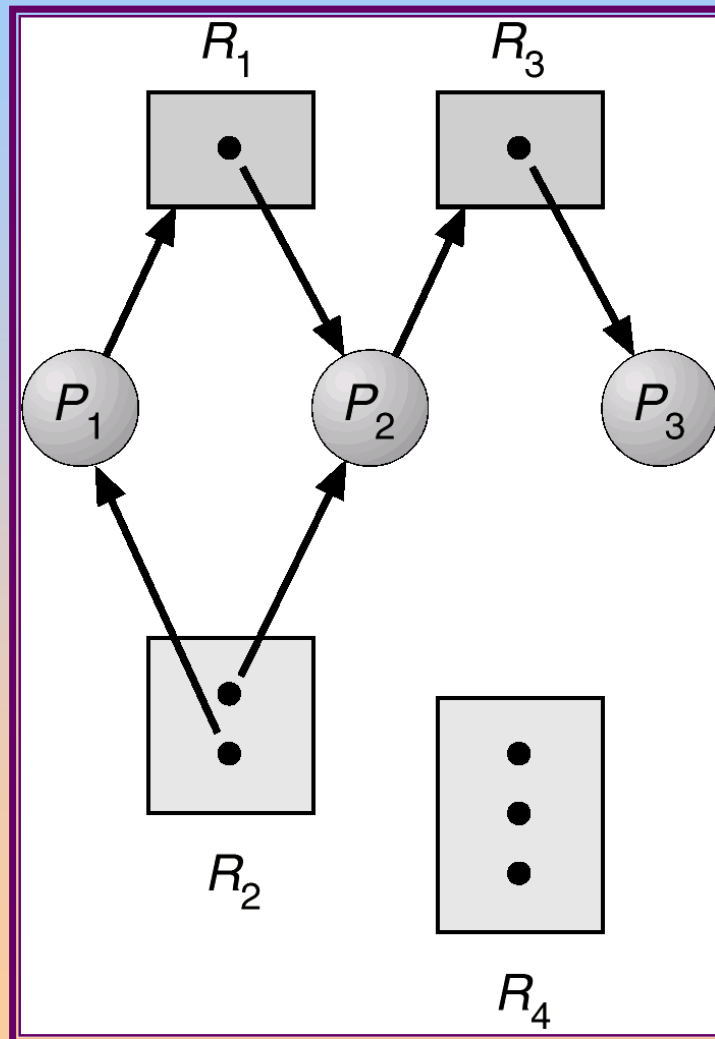
- Process

- Resource Type with 4 instances

- $P_i$ requests instance of $R_j$

    $P_i \longrightarrow \boxed{R_j}$

- $P_i$ is holding an instance of $R_j$

    $P_i \longleftarrow \boxed{R_j}$

# Example of a Resource Allocation Graph
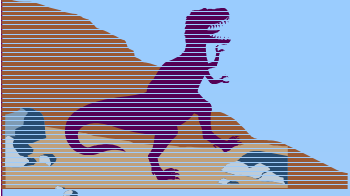
# Resource Allocation Graph With A Deadlock

# Basic Facts

- If graph contains no cycles $\Rightarrow$ no deadlock.

- If graph contains a cycle $\Rightarrow$
  - if only one instance per resource type, then deadlock.
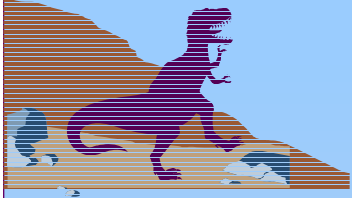  - if several instances per resource type, possibility of deadlock.

# Methods for Handling Deadlocks

- Ensure that the system will *never* enter a deadlock state.

- Allow the system to enter a deadlock state and then recover.

- Ignore the problem and pretend that deadlocks never occur in the system; used by most operating systems, including UNIX.

# Deadlock Prevention

Restrain the ways request can be made.

- **Mutual Exclusion** – not required for sharable resources; must hold for nonsharable resources.

- **Hold and Wait** – must guarantee that whenever a process requests a resource, it does not hold any other resources.
  - ✦ Require process to request and be allocated all its resources before it begins execution, or allow process to request resources only when the process has none.
  - ✦ Low resource utilization; starvation possible.

# Deadlock Prevention (Cont.)

- **No Preemption** –
  - ✦ If a process that is holding some resources requests another resource that cannot be immediately allocated to it, then all resources currently being held are released.
  - ✦ Preempted resources are added to the list of resources for which the process is waiting.
  - ✦ Process will be restarted only when it can regain its old resources, as well as the new ones that it is requesting.

- **Circular Wait** – impose a total ordering of all resource types, and require that each process requests resources in an increasing order of enumeration.

# Deadlock Avoidance

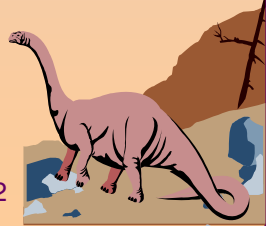Requires that the system has some additional *a priori* information available.

- Simplest and most useful model requires that each process declare the *maximum number* of resources of each type that it may need.

- The deadlock-avoidance algorithm dynamically examines the resource-allocation state to ensure that there can never be a circular-wait condition.

- Resource-allocation *state* is defined by the number of available and allocated resources, and the maximum demands of the processes.

# Safe State

- When a process requests an available resource, system must decide if immediate allocation leaves the system in a *safe state*.

- System is in safe state if there exists a safe sequence of all processes.

- Sequence $<P_1, P_2, …, P_n>$ is safe if for each $P_i$, the resources that $Pi$ can still request can be satisfied by currently available resources + resources held by all the $P_j$, with $j<I$.
  - If $P_i$ resource needs are not immediately available, then $P_i$ can wait until all $P_j$ have finished.
  - When $P_j$ is finished, $P_i$ can obtain needed resources, execute, return allocated resources, and terminate.
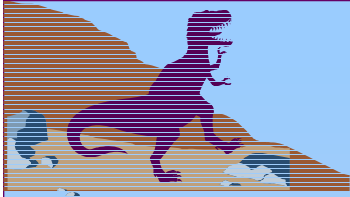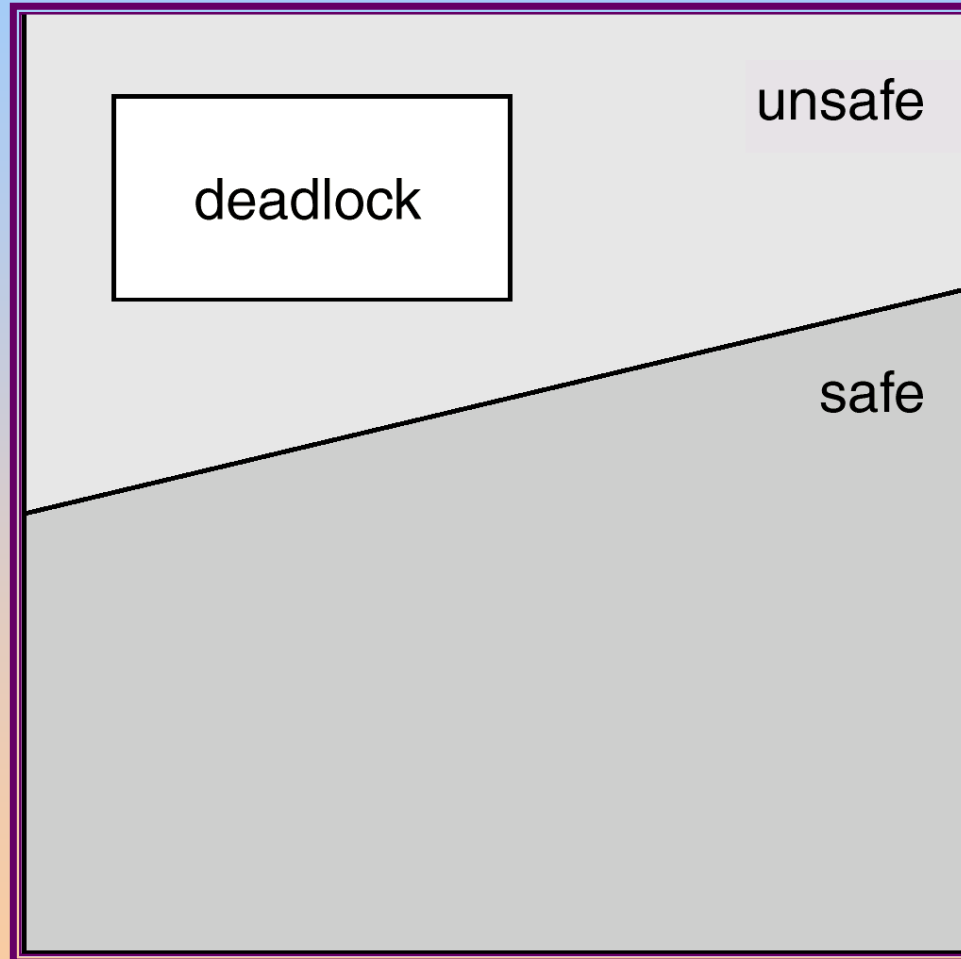  - When $P_i$ terminates, $P_{i+1}$ can obtain its needed resources, and so on.

# Basic Facts

- If a system is in safe state $\Rightarrow$ no deadlocks.

- If a system is in unsafe state $\Rightarrow$ possibility of deadlock.

- Avoidance $\Rightarrow$ ensure that a system will never enter an unsafe state.
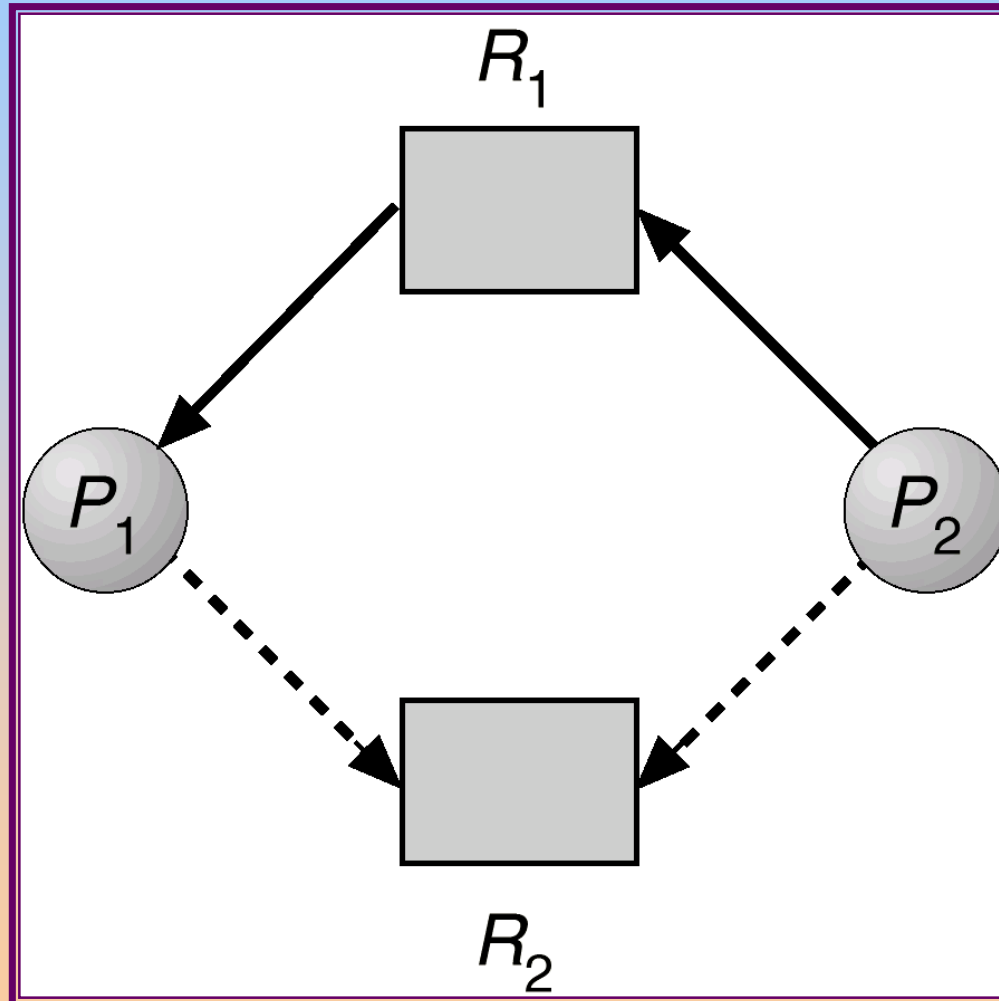
# Safe, Unsafe , Deadlock State
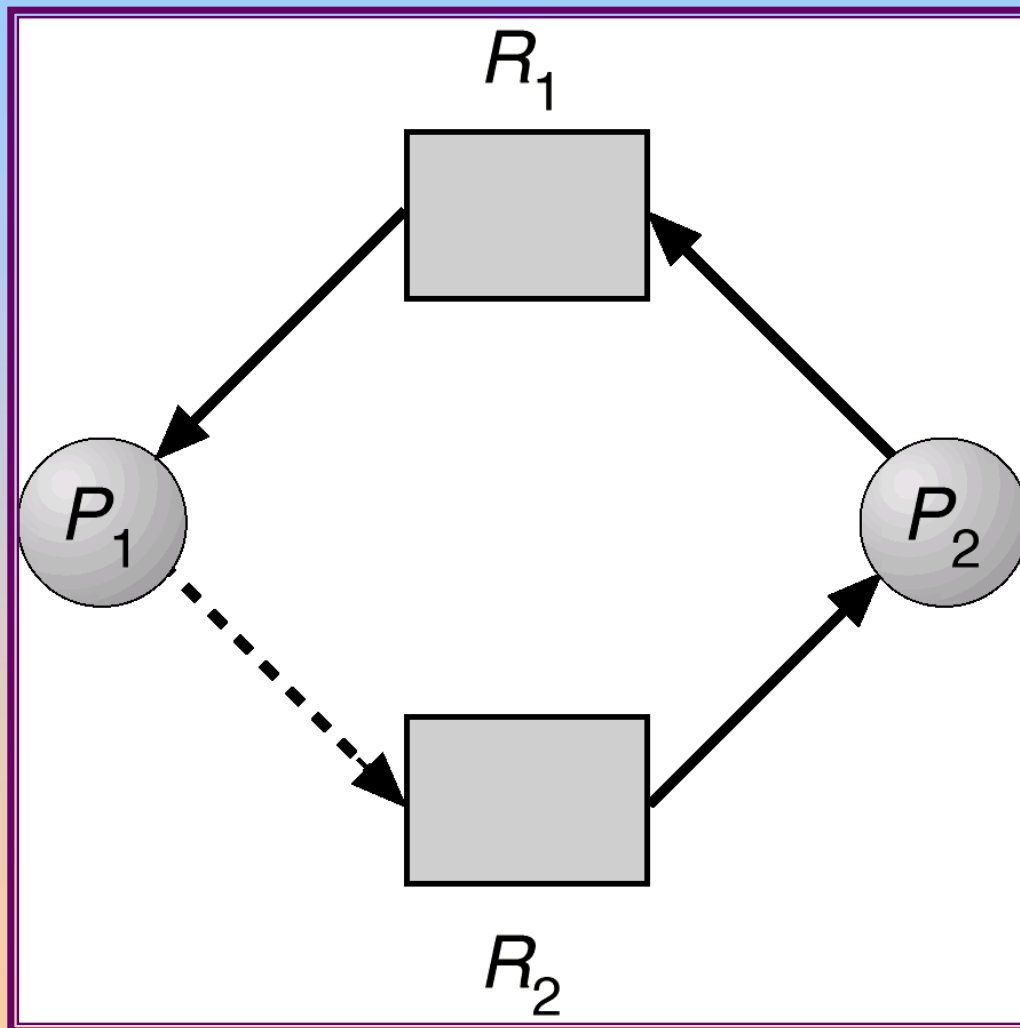
# Resource-Allocation Graph Algorithm

- *Claim edge $P_i \rightarrow R_j$ indicated that process $P_j$ may request resource $R_j$;* represented by a dashed line.

- Claim edge converts to request edge when a process requests a resource.

- When a resource is released by a process, assignment edge reconverts to a claim edge.

- Resources must be claimed *a priori* in the system.

# Unsafe State In Resource-Allocation Graph

# Banker's Algorithm

- Multiple instances.

- Each process must a priori claim maximum use.

- When a process requests a resource it may have to wait.

- When a process gets all its resources it must return them in a finite amount of time.

# Data Structures for the Banker's Algorithm

Let $n$ = number of processes, and $m$ = number of resources types.

- *Available:* Vector of length $m$. If available $[j]$ = $k$, there are $k$ instances of resource type $R_j$ available.

- *Max:* $n$ x $m$ matrix. If *Max* $[i,j]$ = $k$, then process $P_i$ may request at most $k$ instances of resource type $R_j$.

- *Allocation:* $n$ x $m$ matrix. If Allocation$[i,j]$ = $k$ then $P_i$ is currently allocated $k$ instances of $R_j$.

- *Need:* $n$ x $m$ matrix. If *Need*$[i,j]$ = $k$, then $P_i$ may need $k$ more instances of $R_j$ to complete its task.

$$Need\,[i,j] = Max[i,j] - Allocation\,[i,j].$$

# Safety Algorithm

1. Let *Work* and *Finish* be vectors of length *m* and *n*, respectively.  Initialize:

    *Work* = *Available*

    *Finish* [*i*] = *false* for *i* - 1,3, …, *n*.

2. Find and *i* such that both:

    (a) *Finish* [*i*] = *false*

    (b) $Need_i \leq Work$

    If no such *i* exists, go to step 4.

3. *Work* = *Work* + *Allocation_i*
   *Finish*[*i*] = *true*
   go to step 2.

4. If *Finish* [*i*] == true for all *i*, then the system is in a safe state.

# Resource-Request Algorithm for Process $P_i$

*Request* = request vector for process $P_i$.  If $Request_i[j] = k$ then process $P_i$ wants $k$ instances of resource type $R_j$.

1.  If $Request_i \leq Need_i$ go to step 2.  Otherwise, raise error condition, since process has exceeded its maximum claim.

2.  If $Request_i \leq Available$, go to step 3.  Otherwise $P_i$ must wait, since resources are not available.

3.  Pretend to allocate requested resources to $P_i$ by modifying the state as follows:

    $Available = Available = Request_i;$

    $Allocation_i = Allocation_i + Request_i;$

    $Need_i = Need_i - Request_{i;;}$

    -   *If safe $\Rightarrow$ the resources are allocated to $P_i$.*
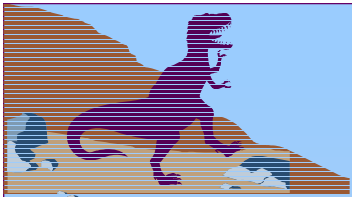    -   *If unsafe $\Rightarrow P_i$ must wait, and the old resource-allocation state is restored*

# Example of Banker's Algorithm

- 5 processes $P_0$ through $P_4$; 3 resource types $A$ (10 instances),
  $B$ (5instances, and $C$ (7 instances).

- Snapshot at time $T_0$:

|  | Allocation | Max | Available |
|---|---|---|---|
|  | A B C | A B C | A B C |
| $P_0$ | 0 1 0 | 7 5 3 | 3 3 2 |
| $P_1$ | 2 0 0 | 3 2 2 |  |
| $P_2$ | 3 0 2 | 9 0 2 |  |
| $P_3$ | 2 1 1 | 2 2 2 |  |
| $P_4$ | 0 0 2 | 4 3 3 |  |

# Example (Cont.)

- The content of the matrix. Need is defined to be Max – Allocation.

$$\underline{Need}$$

|  | A B C |
|---|---|
| $P_0$ | 7 4 3 |
| $P_1$ | 1 2 2 |
| $P_2$ | 6 0 0 |
| $P_3$ | 0 1 1 |
| $P_4$ | 4 3 1 |

- The system is in a safe state since the sequence < $P_1$, $P_3$, $P_4$, $P_2$, $P_0$> satisfies safety criteria.

# Example $P_1$ Request (1,0,2) (Cont.)

■ Check that Request $\leq$ Available (that is, $(1,0,2) \leq (3,3,2) \Rightarrow$ *true.*

|     | Allocation | Need  | Available |
|-----|------------|-------|-----------|
|     | A B C      | A B C | A B C     |
| $P_0$ | 0 1 0    | 7 4 3 | 2 3 0     |
| $P_1$ | 3 0 2    | 0 2 0 |           |
| $P_2$ | 3 0 1    | 6 0 0 |           |
| $P_3$ | 2 1 1    | 0 1 1 |           |
| $P_4$ | 0 0 2    | 4 3 1 |           |

■ Executing safety algorithm shows that sequence $<P_1, P_3, P_4, P_0, P_2>$ satisfies safety requirement.

■ Can request for (3,3,0) by $P_4$ be granted?

■ Can request for (0,2,0) by $P_0$ be granted?

# Deadlock Detection

- Allow system to enter deadlock state

- Detection algorithm

- Recovery scheme

# Single Instance of Each Resource Type

- Maintain *wait-for* graph
  - Nodes are processes.
  - $P_i \rightarrow P_j$ if $P_i$ is waiting for $P_j$.

- Periodically invoke an algorithm that searches for a cycle in the graph.

- An algorithm to detect a cycle in a graph requires an order of $n^2$ operations, where $n$ is the number of vertices in the graph.
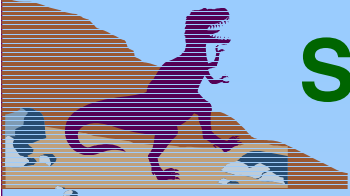
# Resource-Allocation Graph and Wait-for Graph



(a)

(b)

Resource-Allocation Graph          Corresponding wait-for graph

# Several Instances of a Resource Type

- *Available:* A vector of length $m$ indicates the number of available resources of each type.

- *Allocation:* An $n$ x $m$ matrix defines the number of resources of each type currently allocated to each process.

- *Request:* An $n$ x $m$ matrix indicates the current request of each process. If *Request* $[i_j] = k$, then process $P_i$ is requesting $k$ more instances of resource type. $R_j$.

# Detection Algorithm

1. Let *Work* and *Finish* be vectors of length *m* and *n*, respectively Initialize:

   (a) *Work* = *Available*

   (b) For $i$ = 1,2, …, *n*, if *Allocation*$_i$ ≠ 0, then *Finish*[i] = false;otherwise, *Finish*[i] = *true*.

2. Find an index *i* such that both:

   (a) *Finish*[*i*] == *false*

   (b) *Request*$_i$ ≤ *Work*

   If no such *i* exists, go to step 4.

# Detection Algorithm (Cont.)

3. *Work = Work + Allocation$_i$*
   *Finish*[*i*] = *true*
   go to step 2.

4. If *Finish*[*i*] == false, for some *i*, 1 ≤ *i* ≤ *n*, then the system is in deadlock state. Moreover, if *Finish*[*i*] == *false*, then $P_i$ is deadlocked.

Algorithm requires an order of O(*m* x $n^{2)}$ operations to detect whether the system is in deadlocked state.

# Example of Detection Algorithm

- Five processes $P_0$ through $P_4$; three resource types A (7 instances), $B$ (2 instances), and $C$ (6 instances).

- Snapshot at time $T_0$:

|  | Allocation | Request | Available |
|---|---|---|---|
|  | A B C | A B C | A B C |
| $P_0$ | 0 1 0 | 0 0 0 | 0 0 0 |
| $P_1$ | 2 0 0 | 2 0 2 |  |
| $P_2$ | 3 0 3 | 0 0 0 |  |
| $P_3$ | 2 1 1 | 1 0 0 |  |
| $P_4$ | 0 0 2 | 0 0 2 |  |

- Sequence $<P_0, P_2, P_3, P_1, P_4>$ will result in *Finish*[$i$] = true for all $i$.

# Example (Cont.)

- $P_2$ requests an additional instance of type $C$.

  *Request*

  |       | A B C |
  |-------|-------|
  | $P_0$ | 0 0 0 |
  | $P_1$ | 2 0 1 |
  | $P_2$ | 0 0 1 |
  | $P_3$ | 1 0 0 |
  | $P_4$ | 0 0 2 |

- State of system?

  - ✦ Can reclaim resources held by process $P_0$, but insufficient resources to fulfill other processes; requests.
  - ✦ Deadlock exists, consisting of processes $P_1$, $P_2$, $P_3$, and $P_4$.

# Detection-Algorithm Usage

- When, and how often, to invoke depends on:
    - How often a deadlock is likely to occur?
    - How many processes will need to be rolled back?
        - ✔ one for each disjoint cycle

- If detection algorithm is invoked arbitrarily, there may be many cycles in the resource graph and so we would not be able to tell which of the many deadlocked processes "caused" the deadlock.

# Recovery from Deadlock:  Process Termination

- Abort all deadlocked processes.

- Abort one process at a time until the deadlock cycle is eliminated.

- In which order should we choose to abort?
  - Priority of the process.
  - How long process has computed, and how much longer to completion.
  - Resources the process has used.
  - Resources process needs to complete.
  - How many processes will need to be terminated.
  - Is process interactive or batch?

# Recovery from Deadlock: Resource Preemption

- Selecting a victim – minimize cost.

- Rollback – return to some safe state, restart process for that state.

- Starvation – same process may always be picked as victim, include number of rollback in cost factor.

# Combined Approach to Deadlock Handling

■ Combine the three basic approaches

   ✦ prevention

   ✦ avoidance

   ✦ detection

   allowing the use of the optimal approach for each of resources in the system.

■ Partition resources into hierarchically ordered classes.

■ Use most appropriate technique for handling deadlocks within each class.

# Traffic Deadlock for Exercise 8.4

# Chapter 9:  Memory Management

- Background

- Swapping

- Contiguous Allocation

- Paging

- Segmentation

- Segmentation with Paging

# Background

- Program must be brought into memory and placed within a process for it to be run.

- *Input queue* – collection of processes on the disk that are waiting to be brought into memory to run the program.

- User programs go through several steps before being run.

# Binding of Instructions and Data to Memory

Address binding of instructions and data to memory addresses can happen at three different stages.

- **Compile time**:  If memory location known a priori, absolute code can be generated; must recompile code if starting location changes.

- **Load time**:  Must generate *relocatable* code if memory location is not known at compile time.

- **Execution time**:  Binding delayed until run time if the process can be moved during its execution from one memory segment to another.  Need hardware support for address maps (e.g., *base* and *limit registers*).

# Multistep Processing of a User Program

# Logical vs. Physical Address Space

- The concept of a logical *address space* that is bound to a separate *physical address space* is central to proper memory management.
  - *Logical address* – generated by the CPU; also referred to as *virtual address*.
  - *Physical address* – address seen by the memory unit.

- Logical and physical addresses are the same in compile-time and load-time address-binding schemes; logical (virtual) and physical addresses differ in execution-time address-binding scheme.

# Memory-Management Unit (MMU)

- Hardware device that maps virtual to physical address.

- In MMU scheme, the value in the relocation register is added to every address generated by a user process at the time it is sent to memory.

- The user program deals with *logical* addresses; it never sees the *real* physical addresses.

# Dynamic relocation using a relocation register

# Dynamic Loading

- Routine is not loaded until it is called

- Better memory-space utilization; unused routine is never loaded.

- Useful when large amounts of code are needed to handle infrequently occurring cases.

- No special support from the operating system is required implemented through program design.

# Dynamic Linking

- Linking postponed until execution time.

- Small piece of code, *stub*, used to locate the appropriate memory-resident library routine.

- Stub replaces itself with the address of the routine, and executes the routine.

- Operating system needed to check if routine is in processes' memory address.

- Dynamic linking is particularly useful for libraries.

# Overlays

- Keep in memory only those instructions and data that are needed at any given time.

- Needed when process is larger than amount of memory allocated to it.

- Implemented by user, no special support needed from operating system, programming design of overlay structure is complex

# Overlays for a Two-Pass Assembler

# Swapping

- A process can be *swapped* temporarily out of memory to a *backing store*, and then brought back into memory for continued execution.

- Backing store – fast disk large enough to accommodate copies of all memory images for all users; must provide direct access to these memory images.

- *Roll out, roll in* – swapping variant used for priority-based scheduling algorithms; lower-priority process is swapped out so higher-priority process can be loaded and executed.

- Major part of swap time is transfer time; total transfer time is directly proportional to the *amount* of memory swapped.

- Modified versions of swapping are found on many systems, i.e., UNIX, Linux, and Windows.

# Schematic View of Swapping



operating system

① swap out

② swap in

user space

main memory

process $P_1$

process $P_2$

backing store

# Contiguous Allocation

- Main memory usually into two partitions:
  - Resident operating system, usually held in low memory with interrupt vector.
  - User processes then held in high memory.

- Single-partition allocation
  - Relocation-register scheme used to protect user processes from each other, and from changing operating-system code and data.
  - Relocation register contains value of smallest physical address; limit register contains range of logical addresses – each logical address must be less than the limit register.

# Hardware Support for Relocation and Limit Registers

# Contiguous Allocation (Cont.)

- Multiple-partition allocation
  - *Hole* – block of available memory; holes of various size are scattered throughout memory.
  - When a process arrives, it is allocated memory from a hole large enough to accommodate it.
  - Operating system maintains information about:
    a) allocated partitions    b) free partitions (hole)

| OS |
| --- |
| process 5 |
| |
| process 8 |
| |
| process 2 |

→

| OS |
| --- |
| process 5 |
| |
| |
| process 2 |

→

| OS |
| --- |
| process 5 |
| process 9 |
| |
| process 2 |

→

| OS |
| --- |
| process 5 |
| process 9 |
| process 10 |
| |
| process 2 |

# Dynamic Storage-Allocation Problem

How to satisfy a request of size *n* from a list of free holes.

- **First-fit**:  Allocate the *first* hole that is big enough.
- **Best-fit**:  Allocate the *smallest* hole that is big enough; must search entire list, unless ordered by size. Produces the smallest leftover hole.
- **Worst-fit**:  Allocate the *largest* hole; must also search entire list.  Produces the largest leftover hole.

First-fit and best-fit better than worst-fit in terms of speed and storage utilization.

# Fragmentation

- **External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous.

- **Internal Fragmentation** – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used.

- Reduce external fragmentation by compaction
  - Shuffle memory contents to place all free memory together in one large block.
  - Compaction is possible *only* if relocation is dynamic, and is done at execution time.
  - I/O problem
    - Latch job in memory while it is involved in I/O.
    - Do I/O only into OS buffers.

# Paging

- Logical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available.

- Divide physical memory into fixed-sized blocks called **frames** (size is power of 2, between 512 bytes and 8192 bytes).

- Divide logical memory into blocks of same size called **pages**.

- Keep track of all free frames.

- To run a program of size $n$ pages, need to find $n$ free frames and load program.

- Set up a page table to translate logical to physical addresses.

- Internal fragmentation.

# Address Translation Scheme

- Address generated by CPU is divided into:
    - *Page number (p)* – used as an index into a *page table* which contains base address of each page in physical memory.

    - *Page offset (d)* – combined with base address to define the physical memory address that is sent to the memory unit.

# Address Translation Architecture



logical address

physical address

CPU

p | d

f | d

p {

f

page table

f0000 . . . 0000

f1111 . . . 1111

physical memory

f

# Paging Example



page 0
page 1
page 2
page 3

logical memory

page table

| 0 | 1 |
| 1 | 4 |
| 2 | 3 |
| 3 | 7 |

frame number

| 0 | |
| 1 | page 0 |
| 2 | |
| 3 | page 2 |
| 4 | page 1 |
| 5 | |
| 6 | |
| 7 | page 3 |

physical memory

# Paging Example

# Free Frames



(a)

(b)

Before allocation

After allocation

# Implementation of Page Table

- Page table is kept in main memory.

- *Page-table base register (*PTBR) points to the page table.

- *Page-table length register* (PRLR) indicates size of the page table.

- In this scheme every data/instruction access requires two memory accesses.  One for the page table and one for the data/instruction.

- The two memory access problem can be solved by the use of a special fast-lookup hardware cache called *associative memory* or *translation look-aside buffers (TLBs)*

# Associative Memory

- Associative memory – parallel search

| Page # | Frame # |
| --- | --- |
|  |  |
|  |  |
|  |  |
|  |  |

Address translation (A´, A´´)

- ✦ If A´ is in associative register, get frame # out.
- ✦ Otherwise get frame # from page table in memory

# Paging Hardware With TLB

# Effective Access Time

- Associative Lookup = $\varepsilon$ time unit
- Assume memory cycle time is 1 microsecond
- Hit ratio – percentage of times that a page number is found in the associative registers; ration related to number of associative registers.
- Hit ratio = $\alpha$
- Effective Access Time (EAT)

$$EAT = (1 + \varepsilon)\, \alpha + (2 + \varepsilon)(1 - \alpha)$$
$$= 2 + \varepsilon - \alpha$$

# Memory Protection

- Memory protection implemented by associating protection bit with each frame.

- *Valid-invalid* bit attached to each entry in the page table:
    - "valid" indicates that the associated page is in the process' logical address space, and is thus a legal page.
    - "invalid" indicates that the page is not in the process' logical address space.

# Valid (v) or Invalid (i) Bit In A Page Table

# Page Table Structure

- Hierarchical Paging

- Hashed Page Tables

- Inverted Page Tables

# Hierarchical Page Tables

- Break up the logical address space into multiple page tables.

- A simple technique is a two-level page table.

# Two-Level Paging Example

- A logical address (on 32-bit machine with 4K page size) is divided into:
  - a page number consisting of 20 bits.
  - a page offset consisting of 12 bits.
- Since the page table is paged, the page number is further divided into:
  - a 10-bit page number.
  - a 10-bit page offset.
- Thus, a logical address is as follows:

| page number | | page offset |
|:---:|:---:|:---:|
| $p_i$ | $p_2$ | $d$ |
| 10 | 10 | 12 |

where $p_i$ is an index into the outer page table, and $p_2$ is the displacement within the page of the outer page table.

# Two-Level Page-Table Scheme

# Address-Translation Scheme

- Address-translation scheme for a two-level 32-bit paging architecture



logical address

| $p_1$ | $p_2$ | d |

$p_1$ { outer-page table

$p_2$ { page of page table

d {

# Hashed Page Tables

- Common in address spaces > 32 bits.

- The virtual page number is hashed into a page table. This page table contains a chain of elements hashing to the same location.

- Virtual page numbers are compared in this chain searching for a match. If a match is found, the corresponding physical frame is extracted.

# Hashed Page Table

# Inverted Page Table

- One entry for each real page of memory.

- Entry consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page.

- Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs.

- Use hash table to limit the search to one — or at most a few — page-table entries.

# Inverted Page Table Architecture

# Shared Pages

- Shared code
  - One copy of read-only (reentrant) code shared among processes (i.e., text editors, compilers, window systems).
  - Shared code must appear in same location in the logical address space of all processes.

- Private code and data
  - Each process keeps a separate copy of the code and data.
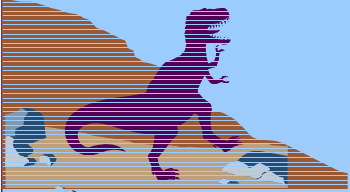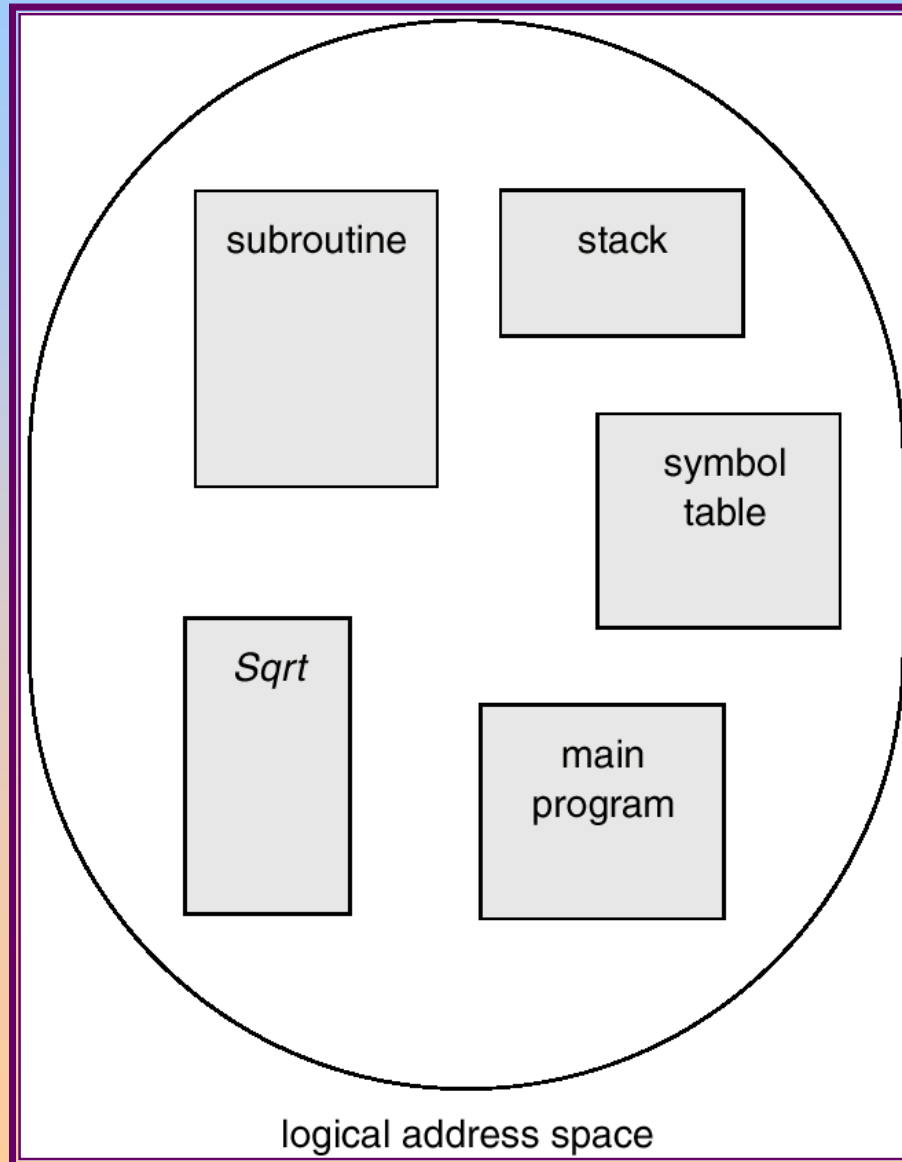  - The pages for the private code and data can appear anywhere in the logical address space.

# Shared Pages Example

| process $P_1$ | | page table for $P_1$ | | | physical memory |
|---|---|---|---|---|---|
| ed 1 | | | 0 | | |
| ed 2 | 3 | | 1 | data 1 | |
| ed 3 | 4 | | 2 | data 3 | |
| data 1 | 6 | | 3 | ed 1 | |
| | 1 | | 4 | ed 2 | |

process $P_1$

page table for $P_1$: 3, 4, 6, 1

process $P_2$: ed 1, ed 2, ed 3, data 2

page table for $P_2$: 3, 4, 6, 7

process $P_3$: ed 1, ed 2, ed 3, data 3

page table for $P_3$: 3, 4, 6, 2

Physical memory:
0
1 data 1
2 data 3
3 ed 1
4 ed 2
5
6 ed 3
7 data 2
8
9
10

# Segmentation

- Memory-management scheme that supports user view of memory.

- A program is a collection of segments.  A segment is a logical unit such as:

  > main program,
  >
  > procedure,
  >
  > function,
  >
  > method,
  >
  > object,
  >
  > local variables, global variables,
  >
  > common block,
  >
  > stack,
  >
  > symbol table, arrays

# User's View of a Program



subroutine

stack

symbol table

Sqrt

main program

logical address space

# Logical View of Segmentation



user space

physical memory space

# Segmentation Architecture

- Logical address consists of a two tuple:

  <segment-number, offset>,

- *Segment table* – maps two-dimensional physical addresses; each table entry has:
  - base – contains the starting physical address where the segments reside in memory.
  - *limit* – specifies the length of the segment.

- *Segment-table base register (STBR)* points to the segment table's location in memory.

- *Segment-table length register (STLR)* indicates number of segments used by a program;

  segment number $s$ is legal if $s <$ STLR.

# Segmentation Architecture (Cont.)

- Relocation.
  - dynamic
  - by segment table

- Sharing.
  - shared segments
  - same segment number

- Allocation.
  - first fit/best fit
  - external fragmentation

# Segmentation Architecture (Cont.)

- Protection.  With each entry in segment table associate:
  - validation bit = 0 $\Rightarrow$ illegal segment
  - read/write/execute privileges
- Protection bits associated with segments; code sharing occurs at segment level.
- Since segments vary in length, memory allocation is a dynamic storage-allocation problem.
- A segmentation example is shown in the following diagram

# Segmentation Hardware

# Example of Segmentation

# Sharing of Segments

# Segmentation with Paging – MULTICS

■ The MULTICS system solved problems of external fragmentation and lengthy search times by paging the segments.

■ Solution differs from pure segmentation in that the segment-table entry contains not the base address of the segment, but rather the base address of a *page table* for this segment.
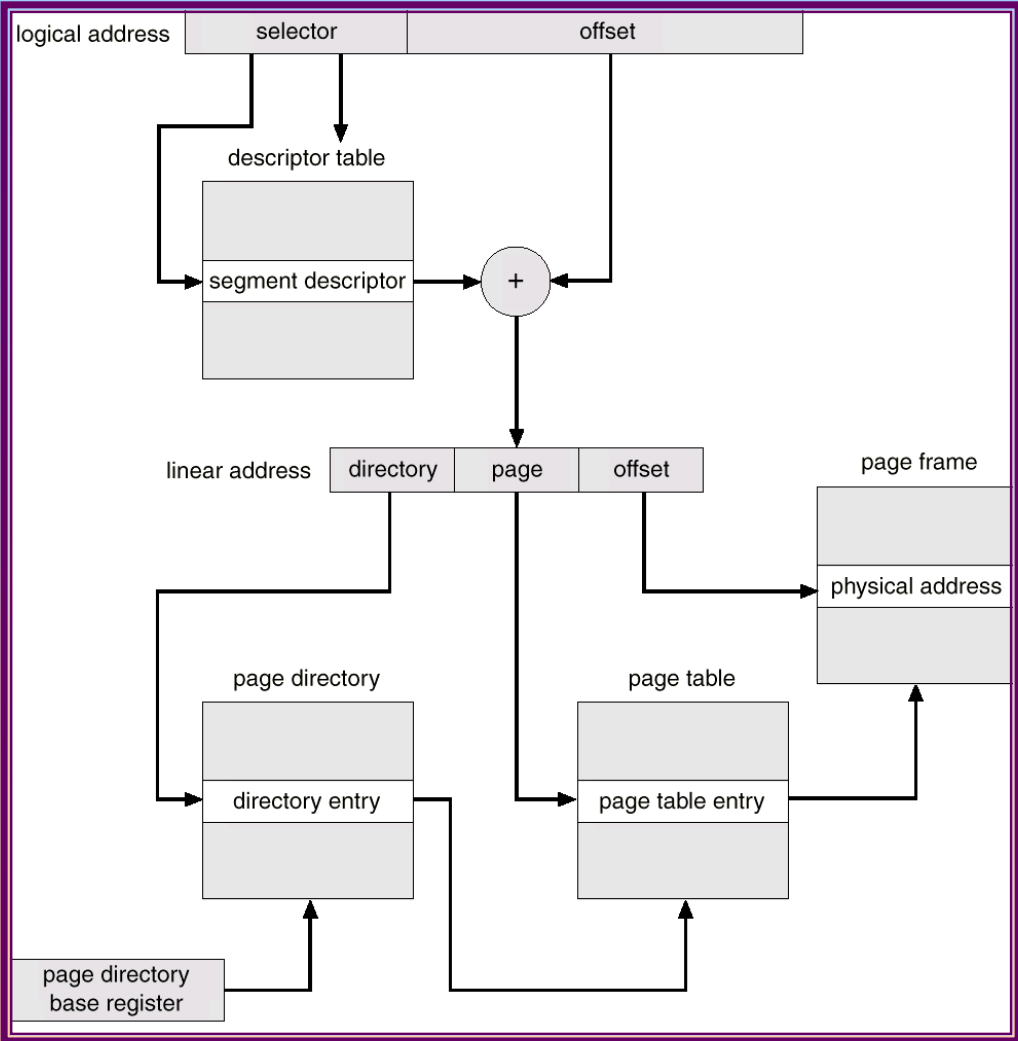
# MULTICS Address Translation Scheme

# **Segmentation with Paging – Intel 386**

■ As shown in the following diagram, the Intel 386 uses segmentation with paging for memory management with a two-level paging scheme.
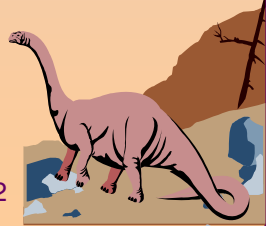
# Intel 30386 Address Translation

# Chapter 10:  Virtual Memory

- Background
- Demand Paging
- Process Creation
- Page Replacement
- Allocation of Frames
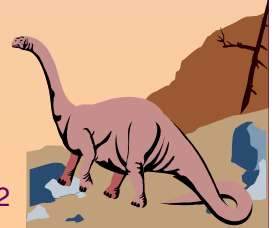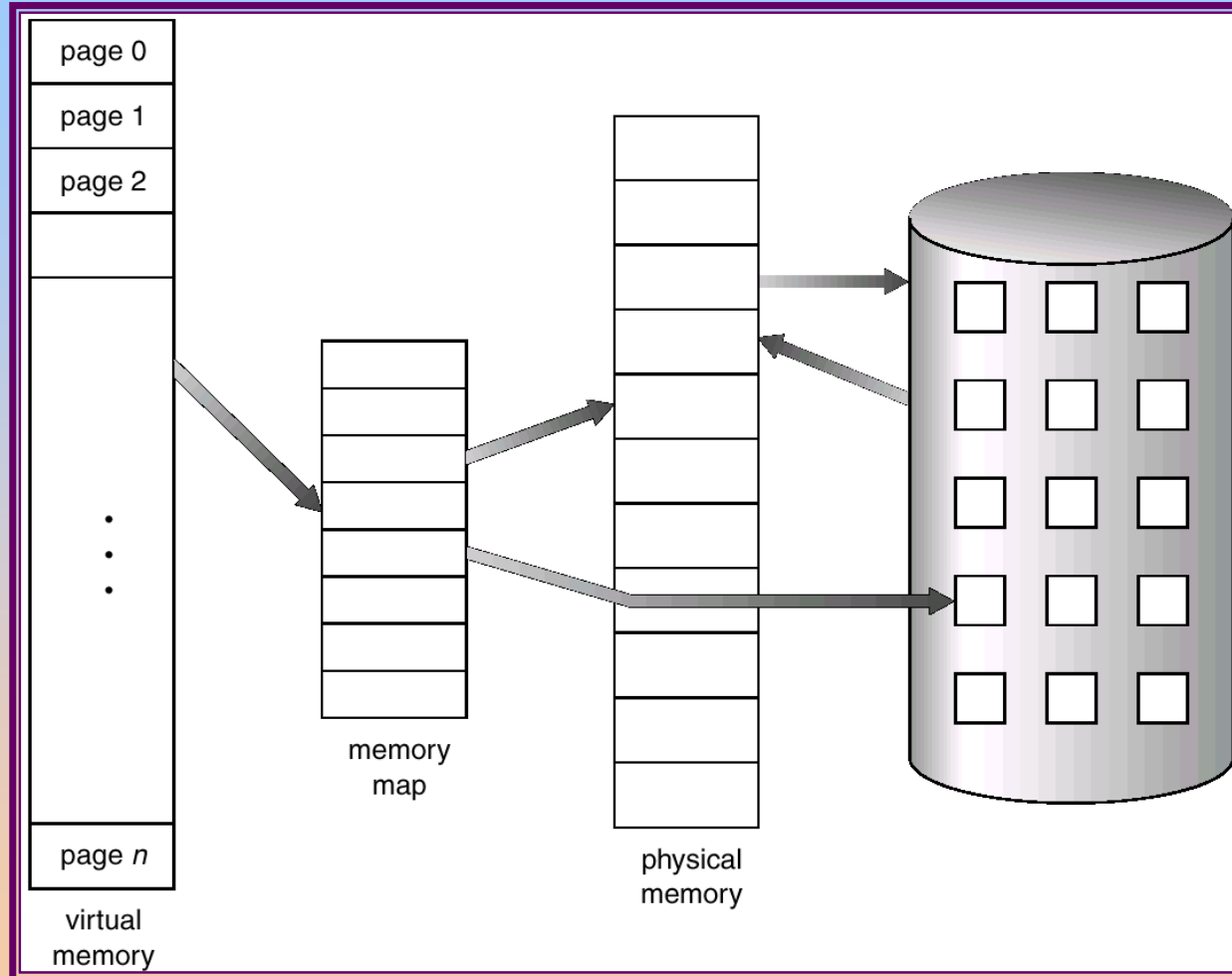- Thrashing
- Operating System Examples

# Background

■ **Virtual memory** – separation of user logical memory from physical memory.

  ✦ Only part of the program needs to be in memory for execution.

  ✦ Logical address space can therefore be much larger than physical address space.

  ✦ Allows address spaces to be shared by several processes.

  ✦ Allows for more efficient process creation.

■ Virtual memory can be implemented via:
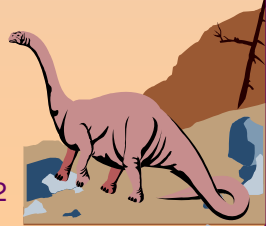
  ✦ Demand paging

  ✦ Demand segmentation

# Virtual Memory That is Larger Than Physical Memory



page 0
page 1
page 2
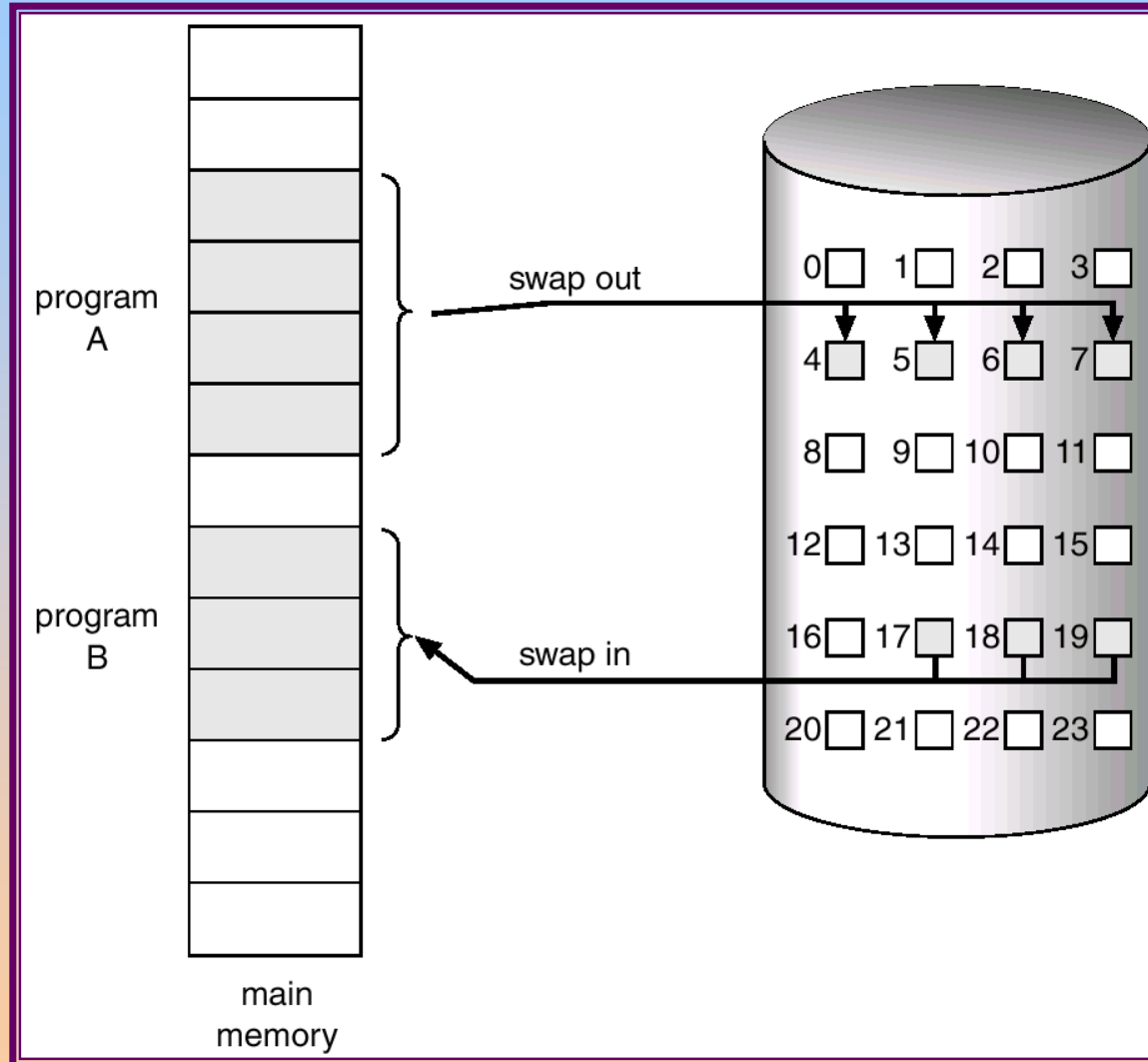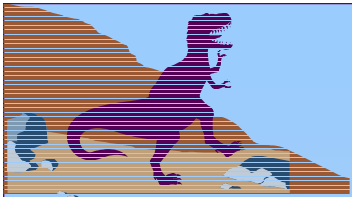
page n

virtual memory

memory map

physical memory

# Demand Paging

- Bring a page into memory only when it is needed.
  - Less I/O needed
  - Less memory needed
  - Faster response
  - More users

- Page is needed $\Rightarrow$ reference to it
  - invalid reference $\Rightarrow$ abort
  - not-in-memory $\Rightarrow$ bring to memory

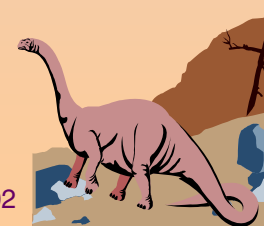# Transfer of a Paged Memory to Contiguous Disk Space

# Valid-Invalid Bit

- With each page table entry a valid–invalid bit is associated
  (1 $\Rightarrow$ in-memory, 0 $\Rightarrow$ not-in-memory)

- Initially valid–invalid but is set to 0 on all entries.

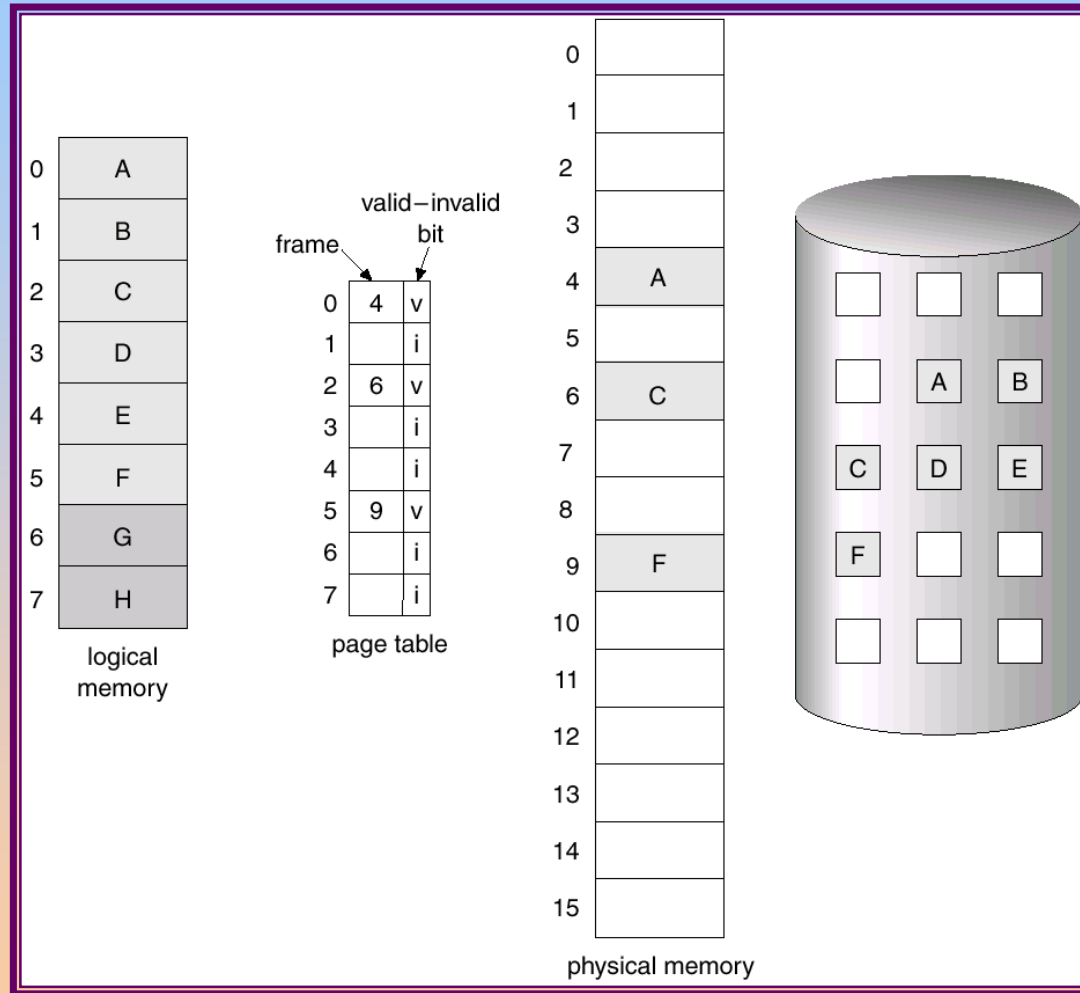- Example of a page table snapshot.

|  Frame #  | valid-invalid bit |
|-----------|:---:|
|           | 1 |
|           | 1 |
|           | 1 |
|           | 1 |
|           | 0 |
| ⋮         |   |
|           | 0 |
|           | 0 |

page table

- During address translation, if valid–invalid bit in page table entry is 0 $\Rightarrow$ page fault.
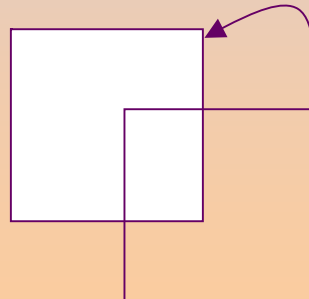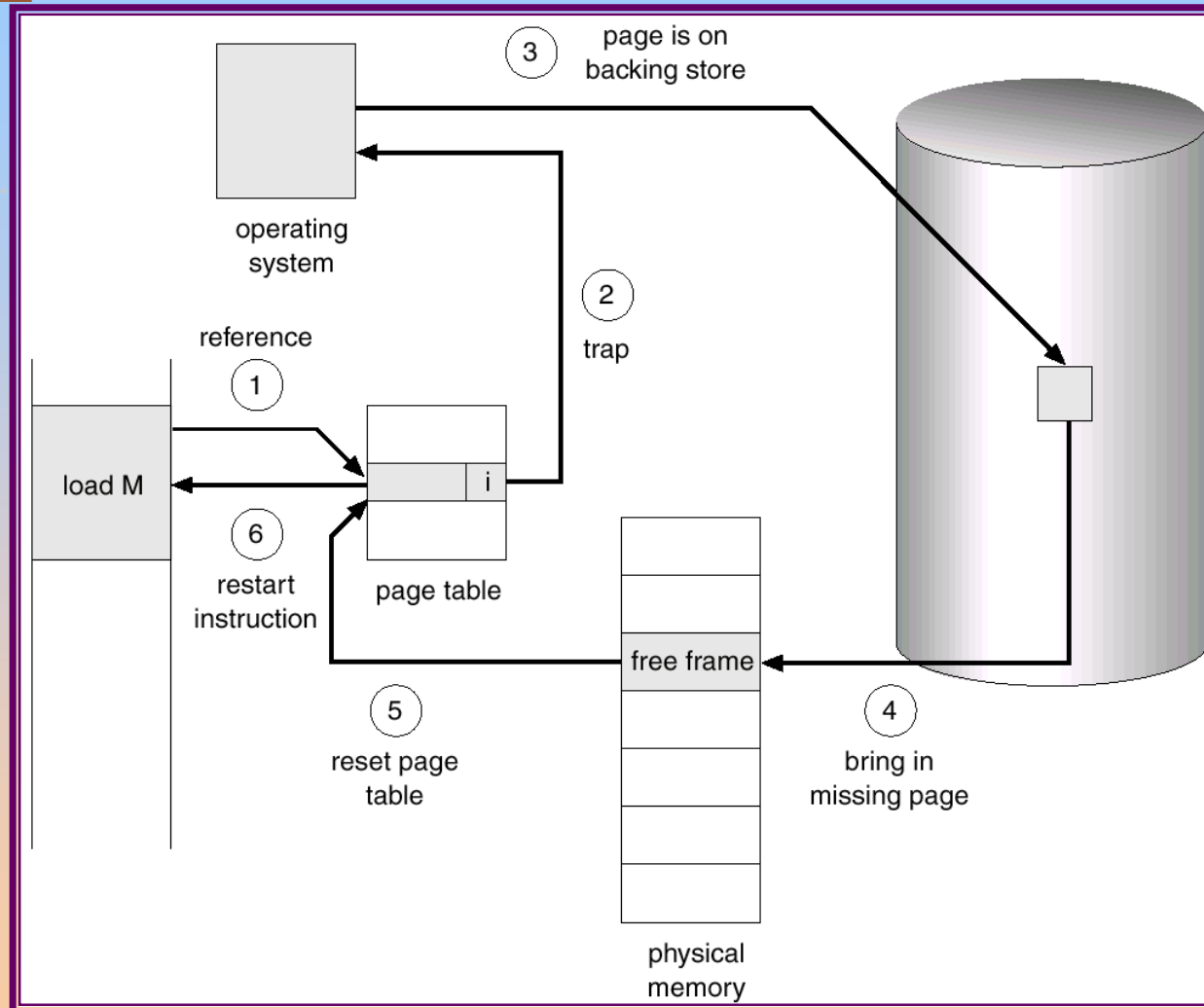
# Page Fault

- If there is ever a reference to a page, first reference will trap to
  OS $\Rightarrow$ page fault

- OS looks at another table to decide:
  - Invalid reference $\Rightarrow$ abort.
  - Just not in memory.

- Get empty frame.

- Swap page into frame.

- Reset tables, validation bit = 1.

- Restart instruction:  Least Recently Used
  - block move

  - auto increment/decrement location

# Steps in Handling a Page Fault



3  page is on backing store

operating system

2  trap

reference

1

load M

6  restart instruction

page table

i

5  reset page table

free frame

4  bring in missing page

physical memory

# What happens if there is no free frame?

- **Page replacement – find some page in memory, but not really in use, swap it out.**
  - algorithm
  - performance – want an algorithm which will result in minimum number of page faults.

- **Same page may be brought into memory several times.**

# Performance of Demand Paging

- Page Fault Rate $0 \leq p \leq 1.0$
  - if $p = 0$ no page faults
  - if $p = 1$, every reference is a fault

- Effective Access Time (EAT)

$$EAT = (1 - p) \text{ x memory access}$$
$$+ p \text{ (page fault overhead}$$
$$+ [\text{swap page out }]$$
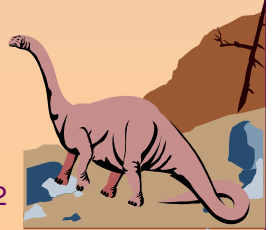$$+ \text{swap page in}$$
$$+ \text{restart overhead)}$$

# Demand Paging Example

- Memory access time = 1 microsecond

- 50% of the time the page that is being replaced has been modified and therefore needs to be swapped out.

- Swap Page Time = 10 msec = 10,000 msec

$$EAT = (1 - p) \times 1 + p \, (15000)$$

$$1 + 15000P \quad \text{(in msec)}$$

# Process Creation

- Virtual memory allows other benefits during process creation:

    - Copy-on-Write

    - Memory-Mapped Files

# Copy-on-Write

- Copy-on-Write (COW) allows both parent and child processes to initially *share* the same pages in memory.

  If either process modifies a shared page, only then is the page copied.

- COW allows more efficient process creation as only modified pages are copied.

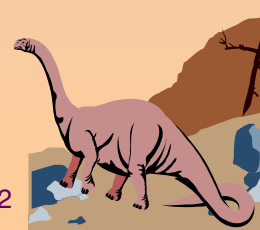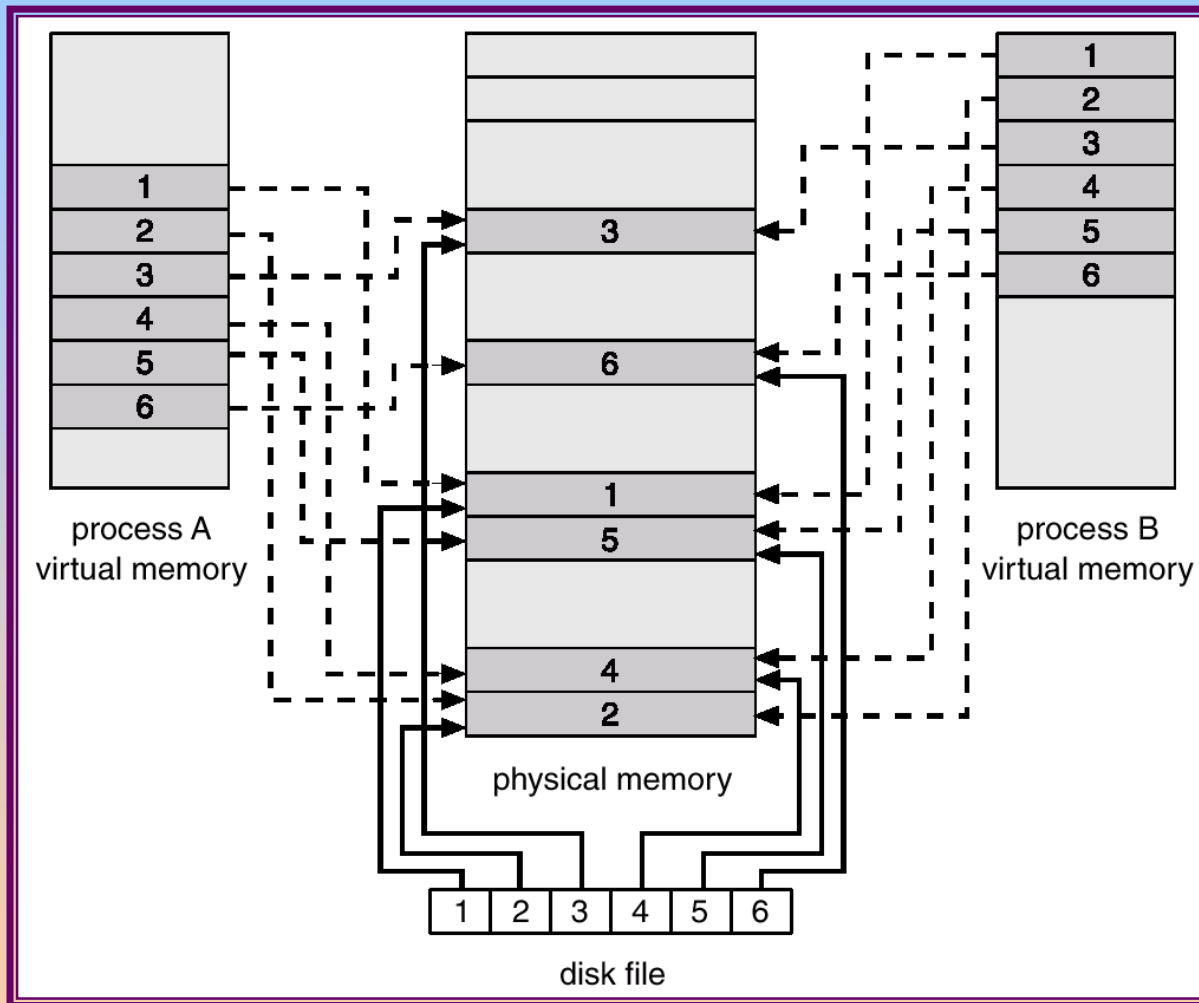- Free pages are allocated from a *pool* of zeroed-out pages.

# Memory-Mapped Files

- Memory-mapped file I/O allows file I/O to be treated as routine memory access by *mapping* a disk block to a page in memory.

- A file is initially read using demand paging. A page-sized portion of the file is read from the file system into a physical page. Subsequent reads/writes to/from the file are treated as ordinary memory accesses.

- Simplifies file access by treating file I/O through memory rather than **read() write()** system calls.

- Also allows several processes to map the same file allowing the pages in memory to be shared.
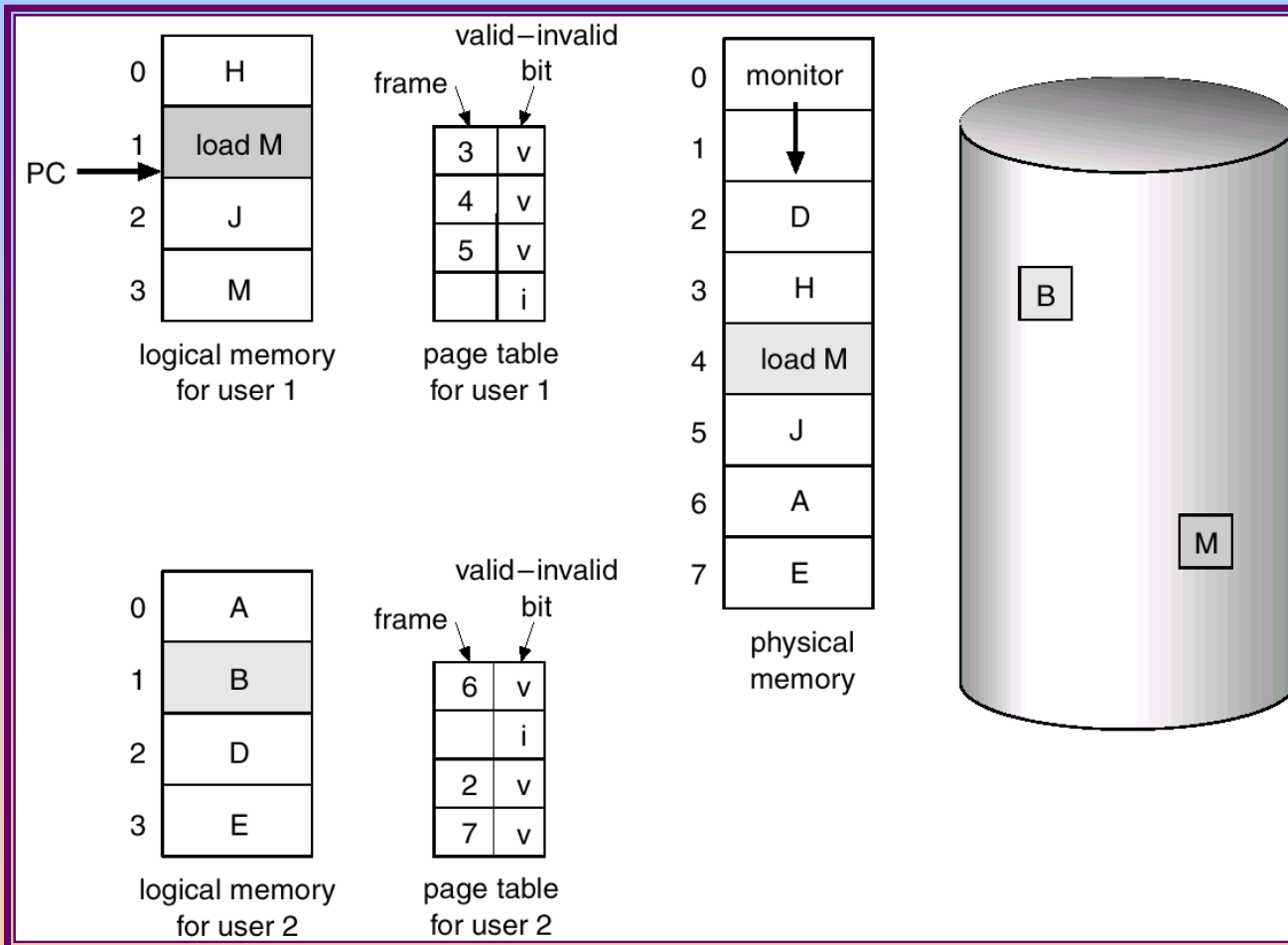
# Memory Mapped Files



process A virtual memory

process B virtual memory

physical memory

disk file

# Page Replacement

- Prevent over-allocation of memory by modifying page-fault service routine to include page replacement.

- Use *modify* (*dirty*) *bit* to reduce overhead of page transfers – only modified pages are written to disk.

- Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory.
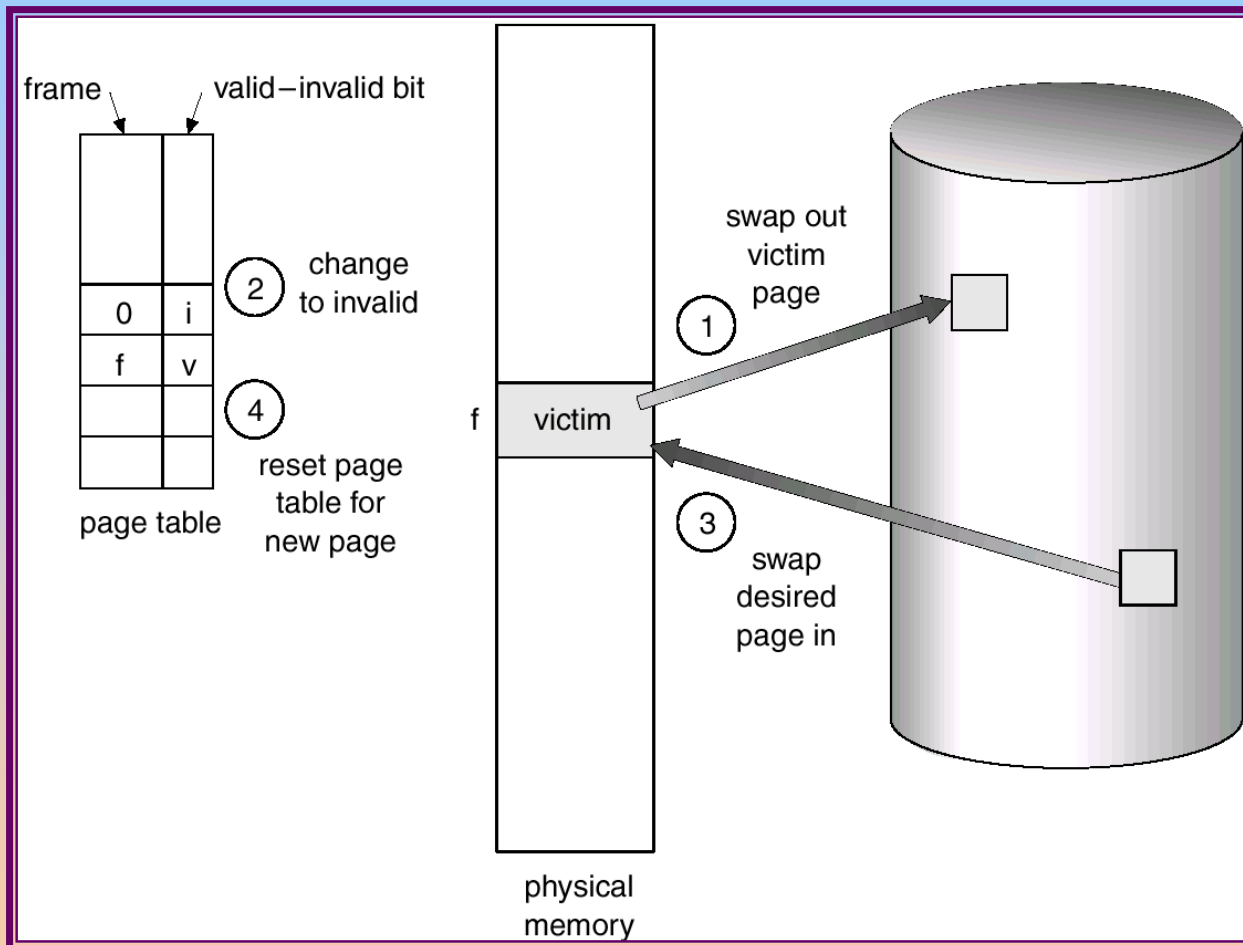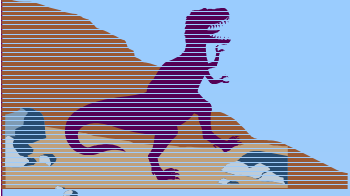
# Need For Page Replacement

# Basic Page Replacement

- Find the location of the desired page on disk.

- Find a free frame:
  - If there is a free frame, use it.
  - If there is no free frame, use a page replacement algorithm to select a *victim* frame.

- Read the desired page into the (newly) free frame. Update the page and frame tables.
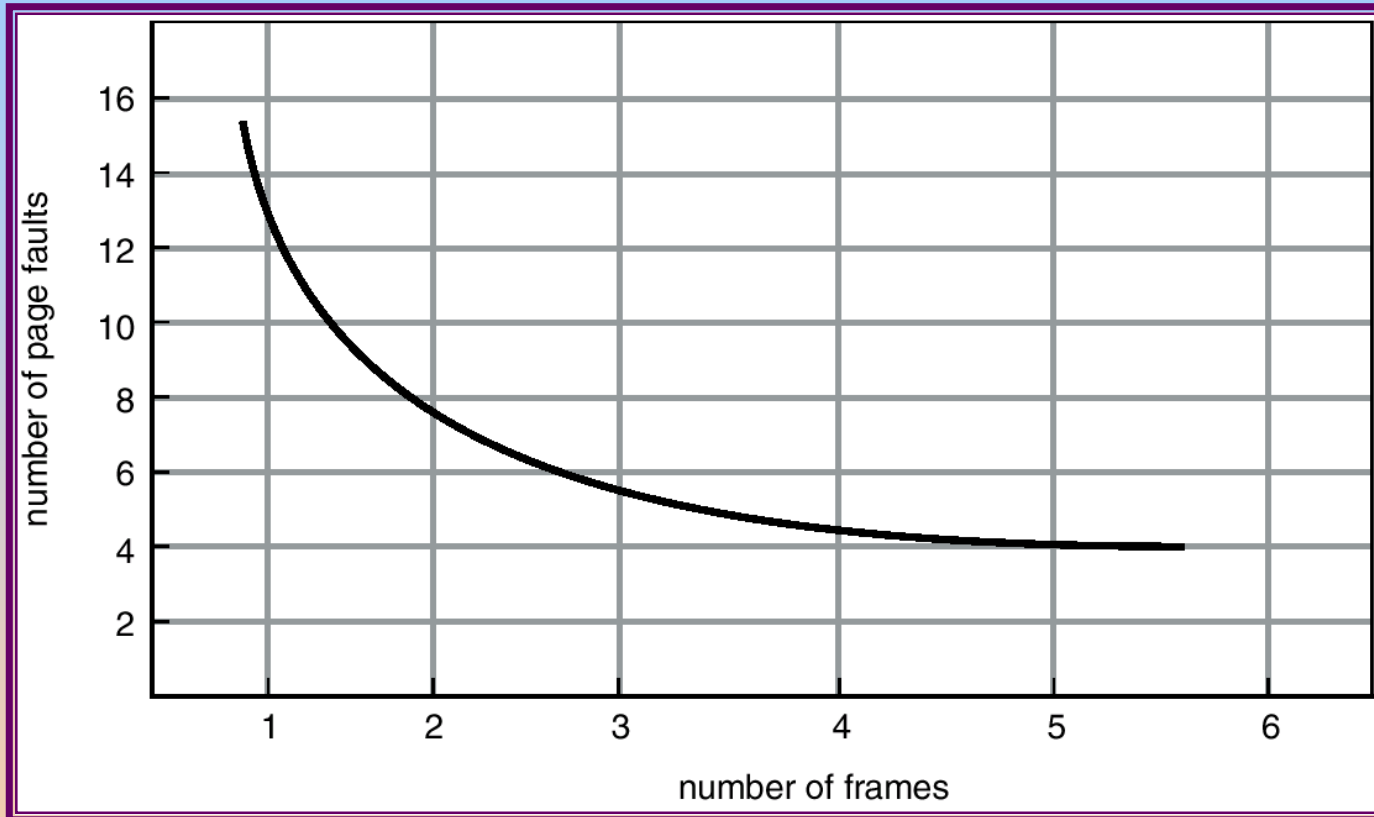
- Restart the process.

# Page Replacement

frame          valid–invalid bit

| 0 | i |
| f | v |
|   |   |
|   |   |

page table

② change to invalid

④ reset page table for new page

f | victim

physical memory

① swap out victim page

③ swap desired page in

# Page Replacement Algorithms

■ Want lowest page-fault rate.

■ Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string.

■ In all our examples, the reference string is

1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5.

# Graph of Page Faults Versus The Number of Frames

# First-In-First-Out (FIFO) Algorithm

- Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5
- 3 frames (3 pages can be in memory at a time per process)

| | | | |
|---|---|---|---|
| 1 | 1 | 4 | 5 |
| 2 | 2 | 1 | 3 |
| 3 | 3 | 2 | 4 |

9 page faults

- 4 frames

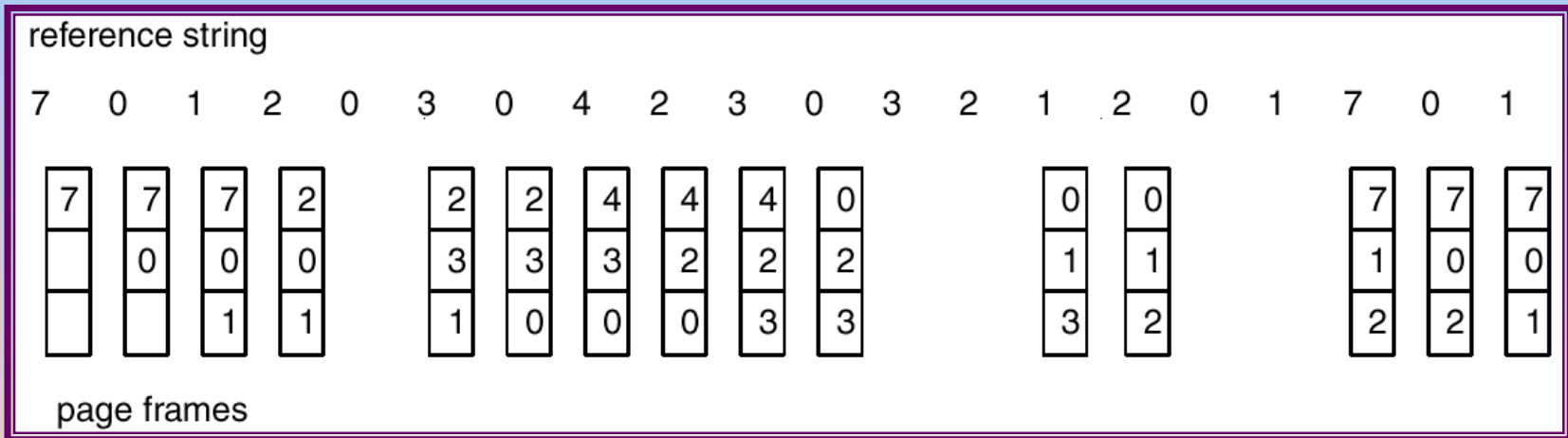| | | | |
|---|---|---|---|
| 1 | 1 | 5 | 4 |
| 2 | 2 | 1 | 5 |
| 3 | 3 | 2 | |
| 4 | 4 | 3 | |

10 page faults
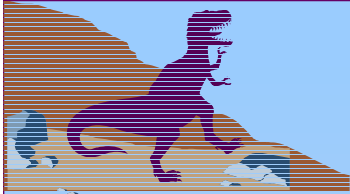
- FIFO Replacement – Belady's Anomaly
  - more frames $\Rightarrow$ less page faults

# FIFO Page Replacement

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

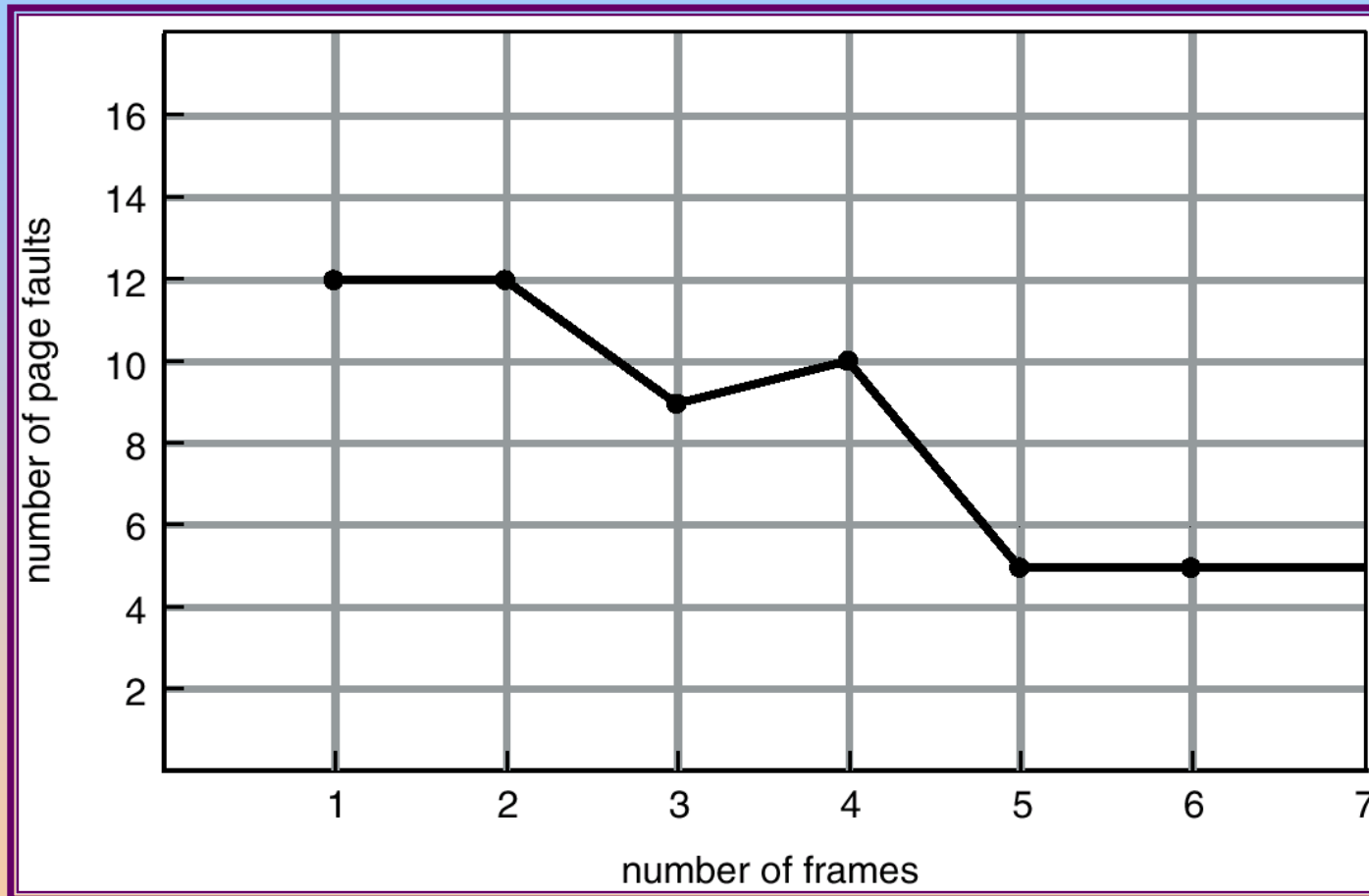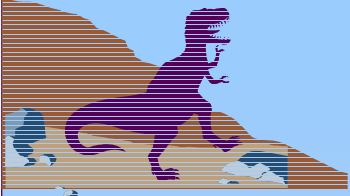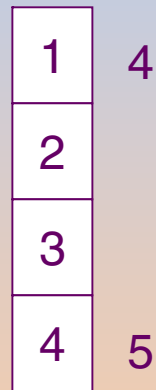| 7 | 7 | 7 | 2 |   | 2 | 2 | 4 | 4 | 4 | 0 |   |   | 0 | 0 |   |   | 7 | 7 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 |   | 3 | 3 | 3 | 2 | 2 | 2 |   |   | 1 | 1 |   |   | 1 | 0 | 0 |
|   |   | 1 | 1 |   | 1 | 0 | 0 | 0 | 3 | 3 |   |   | 3 | 2 |   |   | 2 | 2 | 1 |

page frames

# FIFO Illustrating Belady's Anamoly

# Optimal Algorithm

- Replace page that will not be used for longest period of time.

- 4 frames example

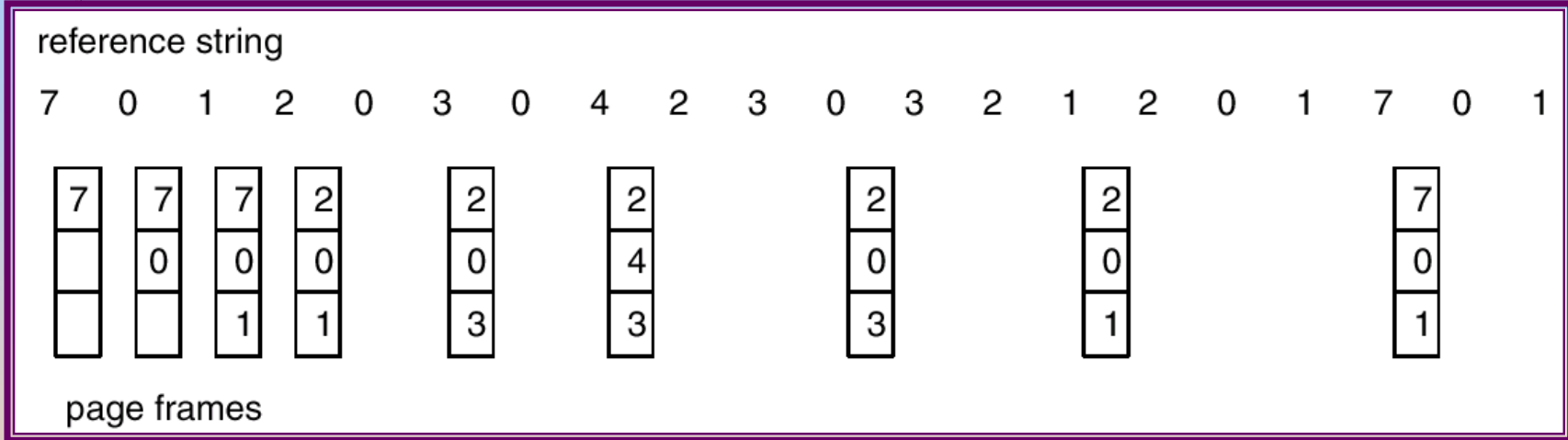  1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

  | 1 | 4 |
  |---|---|
  | 2 | |
  | 3 | |
  | 4 | 5 |

  6 page faults

- How do you know this?

- Used for measuring how well your algorithm performs.
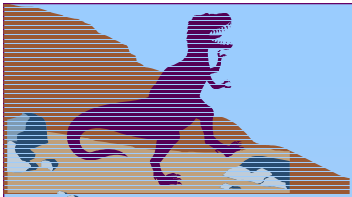
# Optimal Page Replacement

reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

| 7 | 7 | 7 | 2 | | 2 | | 2 | | | 2 | | | 2 | | | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 | | 0 | | 4 | | | 0 | | | 0 | | | 0 |
|   |   | 1 | 1 | | 3 | | 3 | | | 3 | | | 1 | | | 1 |

page frames

# Least Recently Used (LRU) Algorithm

- Reference string:  1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

|   |   |   |
|---|---|---|
| 1 | 5 |   |
| 2 |   |   |
| 3 | 5 | 4 |
| 4 | 3 |   |

- Counter implementation

  - ✦ Every page entry has a counter; every time page is referenced through this entry, copy the clock into the counter.

  - ✦ When a page needs to be changed, look at the counters to determine which are to change.
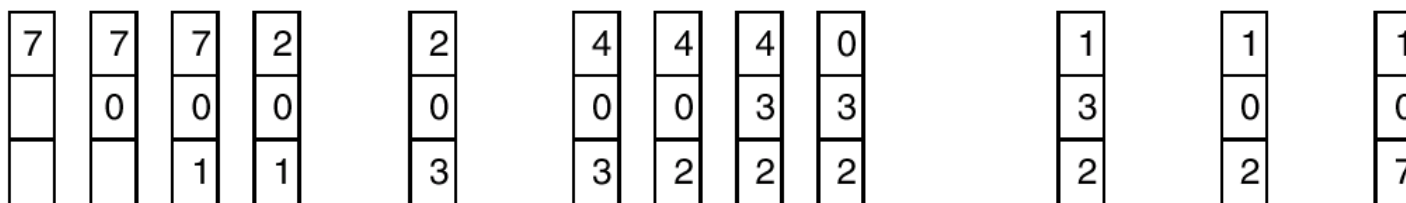
# LRU Page Replacement

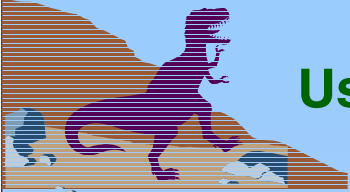reference string

7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

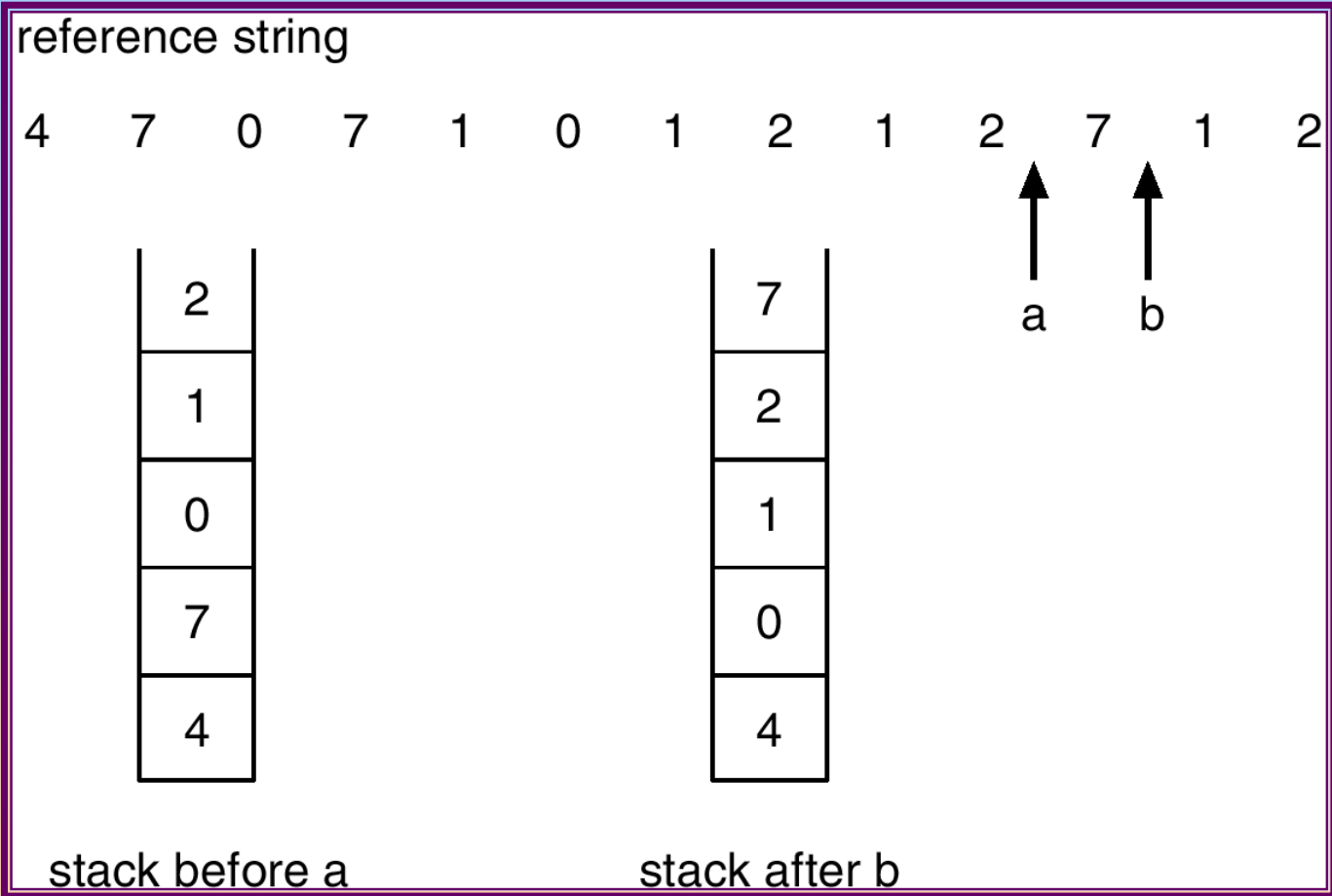| 7 | 7 | 7 | 2 | | 2 | | 4 | 4 | 4 | 0 | | 1 | | 1 | | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 | | 0 | | 0 | 0 | 3 | 3 | | 3 | | 0 | | 0 |
|   |   | 1 | 1 | | 3 | | 3 | 2 | 2 | 2 | | 2 | | 2 | | 7 |

page frames

# LRU Algorithm (Cont.)

■ Stack implementation – keep a stack of page numbers in a double link form:

✦ Page referenced:

✔ move it to the top

✔ requires 6 pointers to be changed

✦ No search for replacement

reference string

4   7   0   7   1   0   1   2   1   2   7   1   2

| 2 |
| 1 |
| 0 |
| 7 |
| 4 |

| 7 |
| 2 |
| 1 |
| 0 |
| 4 |

a       b

stack before a                    stack after b
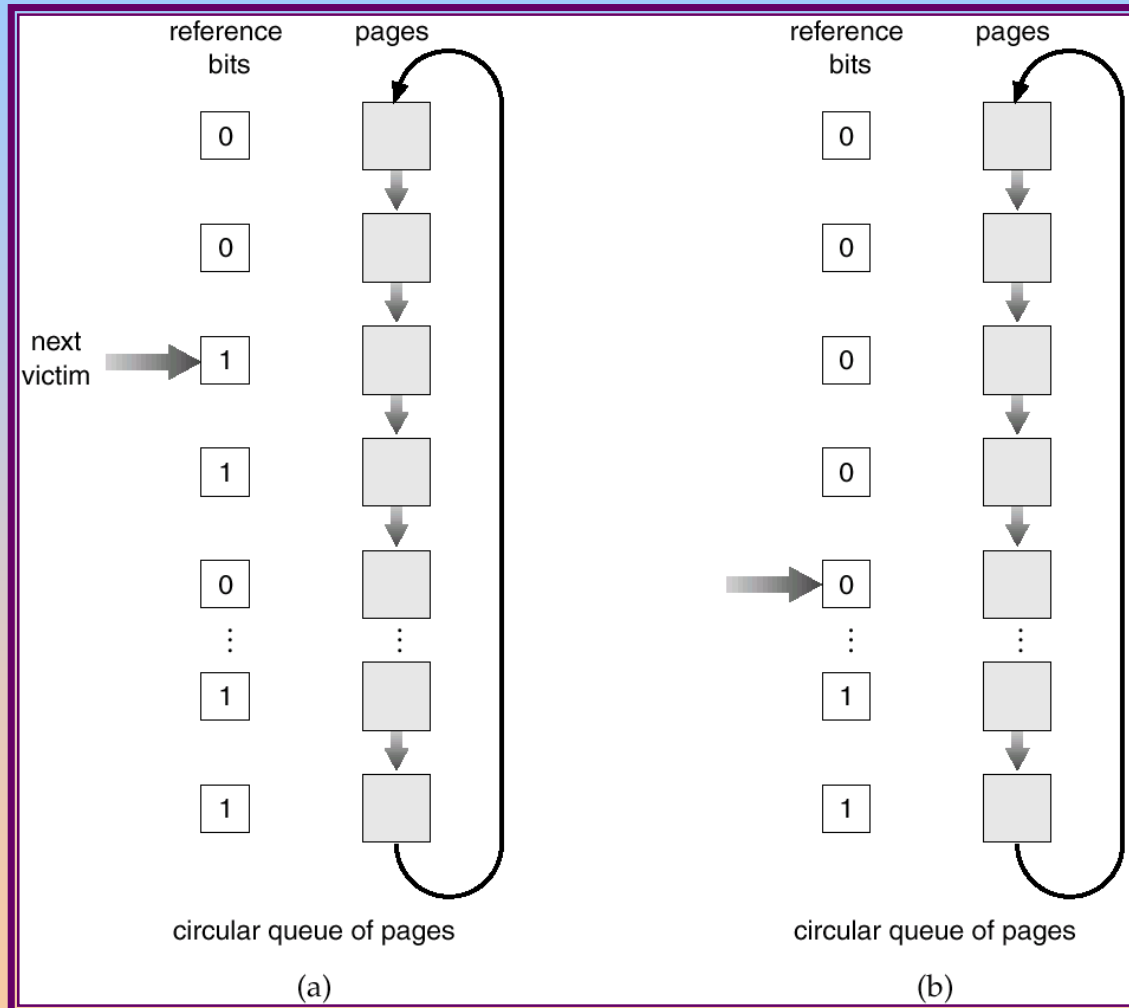
# LRU Approximation Algorithms

- **Reference bit**
  - With each page associate a bit, initially = 0
  - When page is referenced bit set to 1.
  - Replace the one which is 0 (if one exists). We do not know the order, however.

- **Second chance**
  - Need reference bit.
  - Clock replacement.
  - If page to be replaced (in clock order) has reference bit = 1. then:
    - ✔ set reference bit 0.
    - ✔ leave page in memory.
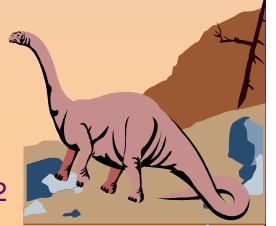    - ✔ replace next page (in clock order), subject to same rules.

# Second-Chance (clock) Page-Replacement Algorithm

reference bits    pages        reference bits    pages

next victim

0

0

1

1

0

⋮

1

1

0

0

0

0

0

⋮

1

1

circular queue of pages

(a)

circular queue of pages

(b)

# Counting Algorithms

■ Keep a counter of the number of references that have been made to each page.

■ LFU Algorithm: replaces page with smallest count.

■ MFU Algorithm: based on the argument that the page with the smallest count was probably just brought in and has yet to be used.

# Allocation of Frames

- Each process needs **minimum** number of pages.
- Example: IBM 370 – 6 pages to handle SS MOVE instruction:
  - instruction is 6 bytes, might span 2 pages.
  - 2 pages to handle **from**.
  - 2 pages to handle **to**.
- Two major allocation schemes.
  - fixed allocation
  - priority allocation

# Fixed Allocation

■ Equal allocation – e.g., if 100 frames and 5 processes, give each 20 pages.

■ Proportional allocation – Allocate according to the size of process.

   – $s_i$ = size of process $p_i$

   – $S = \sum s_i$

   – $m$ = total number of frames

   – $a_i$ = allocation for $p_i = \dfrac{s_i}{S} \times m$
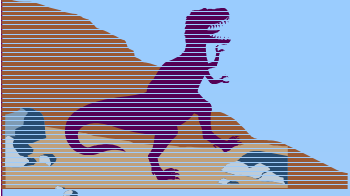
$$m = 64$$

$$s_i = 10$$

$$s_2 = 127$$

$$a_1 = \frac{10}{137} \times 64 \approx 5$$
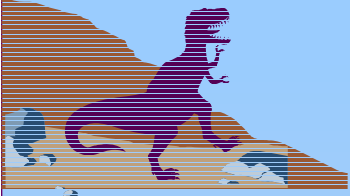
$$a_2 = \frac{127}{137} \times 64 \approx 59$$

# Priority Allocation

■ Use a proportional allocation scheme using priorities rather than size.

■ If process $P_i$ generates a page fault,

   ✦ select for replacement one of its frames.

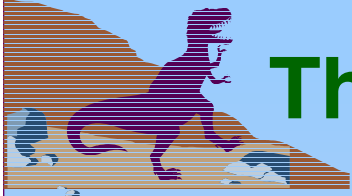   ✦ select for replacement a frame from a process with lower priority number.

# Global vs. Local Allocation

- **Global** replacement – process selects a replacement frame from the set of all frames; one process can take a frame from another.

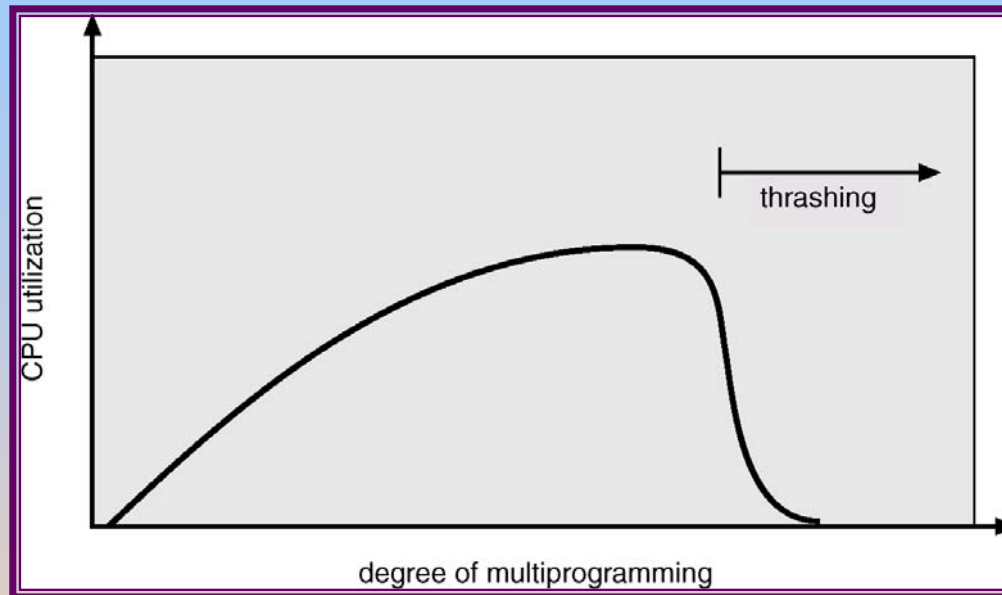- **Local** replacement – each process selects from only its own set of allocated frames.

# Thrashing

■ If a process does not have "enough" pages, the page-fault rate is very high.  This leads to:

 ✦ low CPU utilization.

 ✦ operating system thinks that it needs to increase the degree of multiprogramming.

 ✦ another process added to the system.

■ **Thrashing** $\equiv$ a process is busy swapping pages in and out.
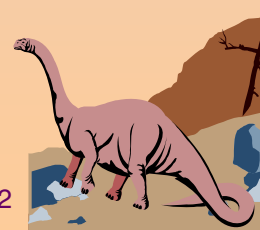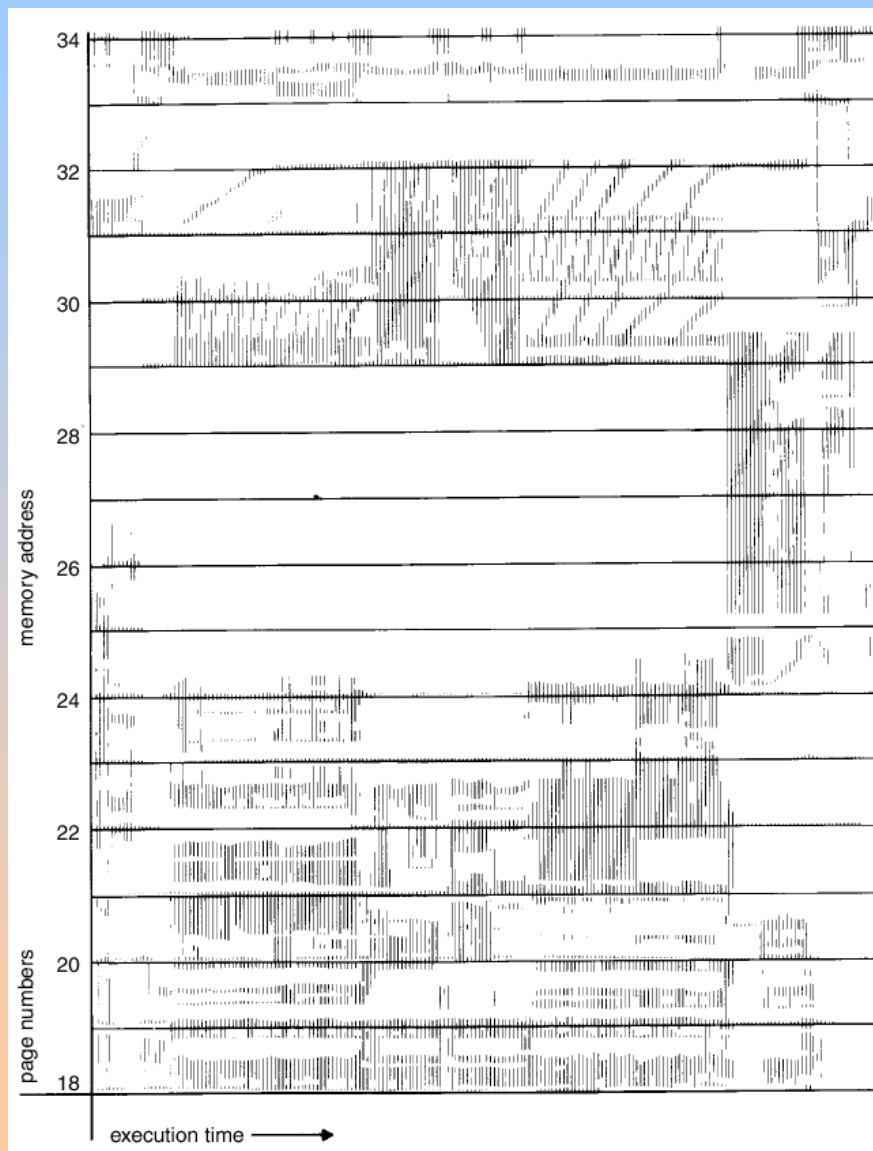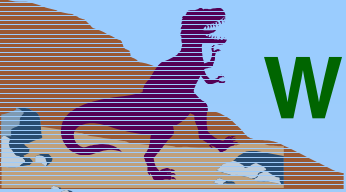
# Thrashing



- Why does paging work?
  Locality model
  - ✦ Process migrates from one locality to another.
  - ✦ Localities may overlap.
- Why does thrashing occur?
  $\Sigma$ size of locality > total memory size

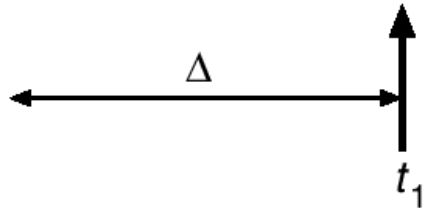# Locality In A Memory-Reference Pattern

# Working-Set Model

- $\Delta \equiv$ working-set window $\equiv$ a fixed number of page references
  Example: 10,000 instruction

- $WSS_i$ (working set of Process $P_i$) =
  total number of pages referenced in the most recent $\Delta$ (varies in time)

  - if $\Delta$ too small will not encompass entire locality.
  - if $\Delta$ too large will encompass several localities.
  - if $\Delta = \infty \Rightarrow$ will encompass entire program.

- $D = \Sigma\ WSS_i \equiv$ total demand frames

- if $D > m \Rightarrow$ Thrashing

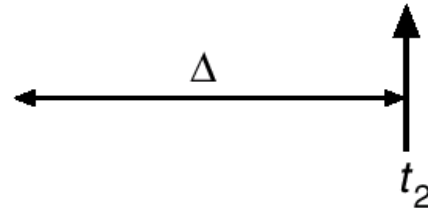- Policy if $D > m$, then suspend one of the processes.

# Working-set model

page reference table

... 2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 ...

$\Delta$

$\Delta$

$t_1$
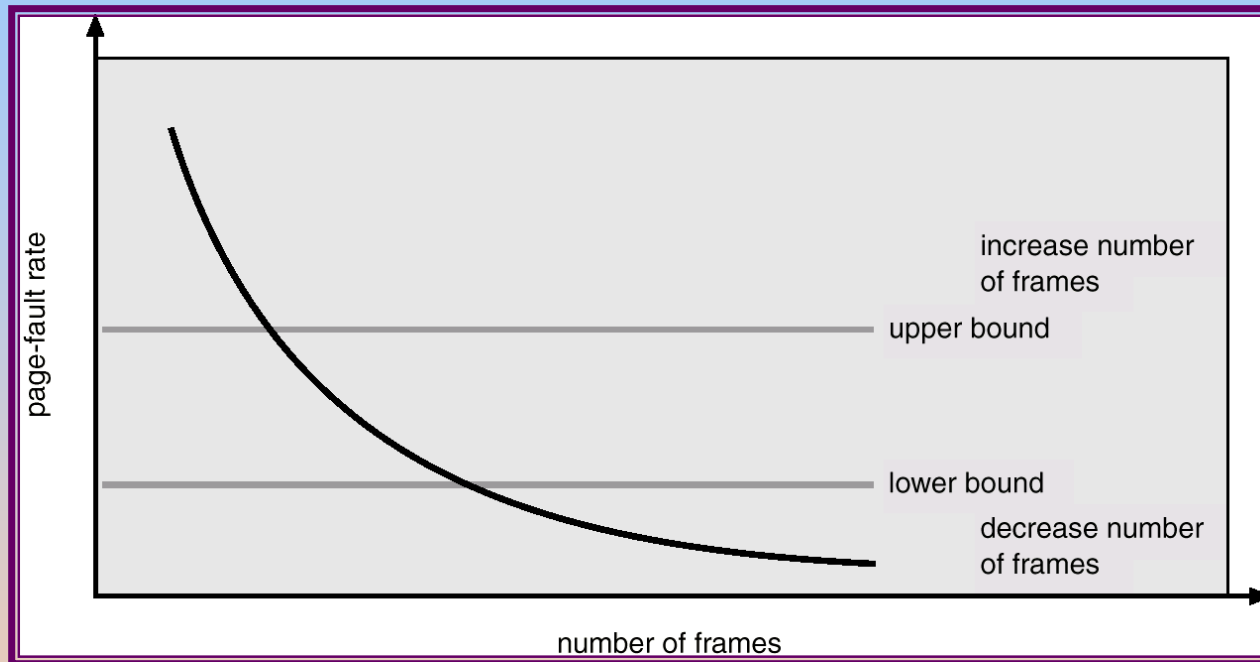
$t_2$

$WS(t_1) = \{1,2,5,6,7\}$

$WS(t_2) = \{3,4\}$

# Keeping Track of the Working Set

- Approximate with interval timer + a reference bit
- Example: $\Delta$ = 10,000
  - Timer interrupts after every 5000 time units.
  - Keep in memory 2 bits for each page.
  - Whenever a timer interrupts copy and sets the values of all reference bits to 0.
  - If one of the bits in memory = 1 $\Rightarrow$ page in working set.
- Why is this not completely accurate?
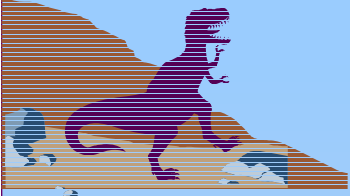- Improvement = 10 bits and interrupt every 1000 time units.
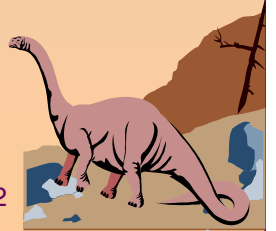
# Page-Fault Frequency Scheme



- Establish "acceptable" page-fault rate.
  - ✦ If actual rate too low, process loses frame.
  - ✦ If actual rate too high, process gains frame.

# Other Considerations

- Prepaging

- Page size selection
  - fragmentation
  - table size
  - I/O overhead
  - locality

# Other Considerations (Cont.)

- **TLB Reach** - The amount of memory accessible from the TLB.

- TLB Reach = (TLB Size) X (Page Size)

- Ideally, the working set of each process is stored in the TLB. Otherwise there is a high degree of page faults.

# Increasing the Size of the TLB

- **Increase the Page Size**. This may lead to an increase in fragmentation as not all applications require a large page size.

- **Provide Multiple Page Sizes**. This allows applications that require larger page sizes the opportunity to use them without an increase in fragmentation.

# Other Considerations (Cont.)

- Program structure

  - **int A[][] = new int[1024][1024];**

  - Each row is stored in one page

  - Program 1
    ```
    for (j = 0; j < A.length; j++)
        for (i = 0; i < A.length; i++)
            A[i,j] = 0;
    ```

    1024 x 1024 page faults

  - Program 2
    ```
    for (i = 0; i < A.length; i++)
        for (j = 0; j < A.length; j++)
            A[i,j] = 0;
    ```
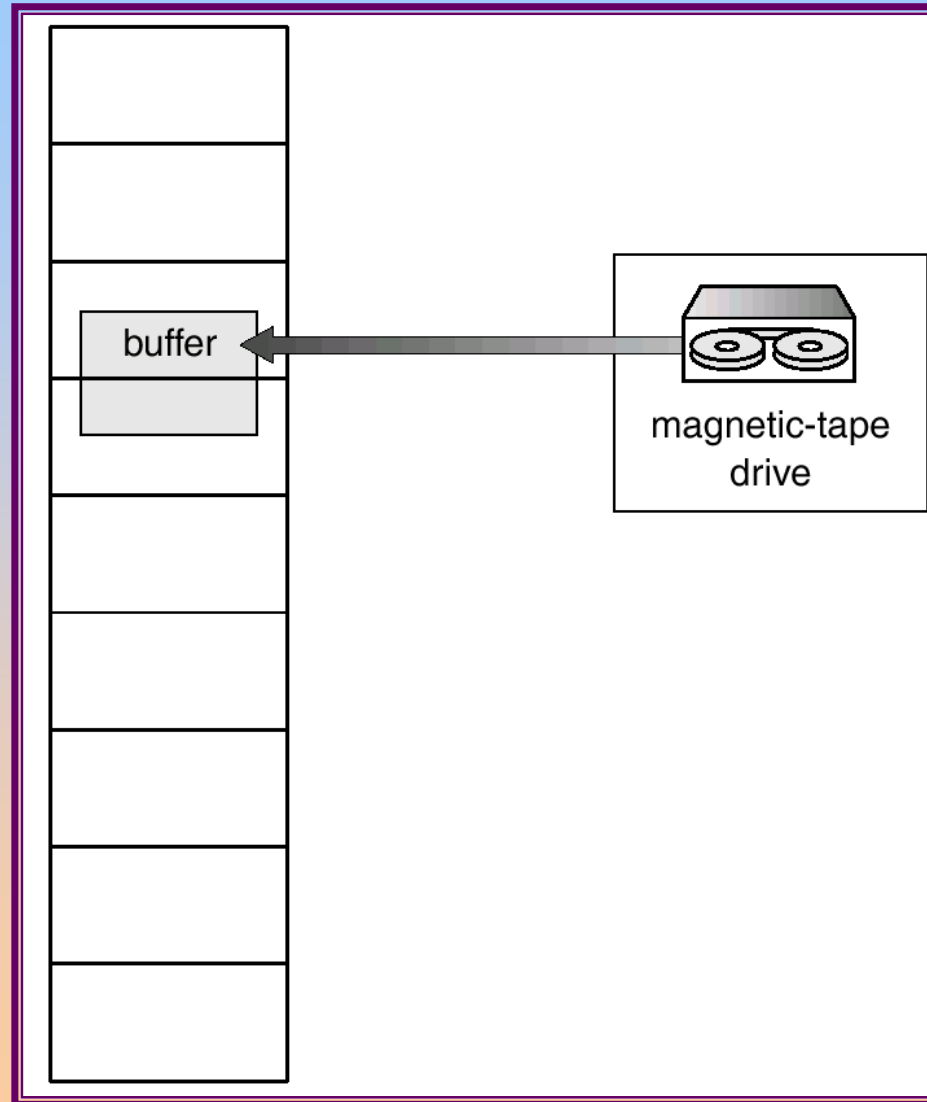
    1024 page faults

# Other Considerations (Cont.)

- **I/O Interlock** – Pages must sometimes be locked into memory.

- Consider I/O. Pages that are used for copying a file from a device must be locked from being selected for eviction by a page replacement algorithm.

buffer

magnetic-tape drive

# Operating System Examples

- Windows NT

- Solaris 2

# Windows NT

- Uses demand paging with **clustering.** Clustering brings in pages surrounding the faulting page.

- Processes are assigned **working set minimum** and **working set maximum**.

- Working set minimum is the minimum number of pages the process is guaranteed to have in memory.

- A process may be assigned as many pages up to its working set maximum.

- When the amount of free memory in the system falls below a threshold, **automatic working set trimming** is performed to restore the amount of free memory.

- Working set trimming removes pages from processes that have pages in excess of their working set minimum.
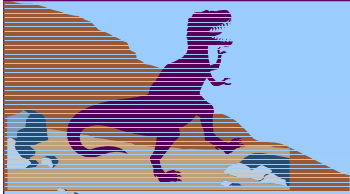
# Solaris 2

- Maintains a list of free pages to assign faulting processes.

- **Lotsfree** – threshold parameter to begin paging.

- Paging is peformed by *pageout* process.

- Pageout scans pages using modified clock algorithm.

- **Scanrate** is the rate at which pages are scanned. This ranged from **slowscan** to **fastscan**.

- Pageout is called more frequently depending upon the amount of free memory available.

# Solar Page Scanner



8192
fastscan

scan rate

100
slowscan

minfree        desfree              lotsfree

amount of free memory