Module 1: Introduction

- What Operating Systems Do
- Computer-System Organization
- Computer-System Architecture
- Operating-System Structure
- Operating-System Operations
- Process Management
- Memory Management
- Storage Management
- Protection and Security
- Distributed Systems
- Special-Purpose Systems
- Computing Environments

Objectives

- To provide a grand tour of the major operating systems components
- To provide coverage of basic computer system organization

What is an Operating System?

- A program that acts as an intermediary between a user of a computer and the computer hardware.
- Operating system goals:
 - Execute user programs and make solving user problems easier.
 - Make the computer system convenient to use.
- Use the computer hardware in an efficient manner.

Computer System Structure

- Computer system can be divided into four components
 - Hardware provides basic computing resources
 - CPU, memory, I/O devices
 - Operating system
 - Controls and coordinates use of hardware among various applications and users
 - Application programs define the ways in which the system resources are used to solve the computing problems of the users
 - Word processors, compilers, web browsers, database systems, video games
 - Users
 - People, machines, other computers

Four Components of a Computer System



Operating System Definition

- OS is a **resource allocator**
 - Manages all resources
 - Decides between conflicting requests for efficient and fair resource use
- OS is a control program
 - Controls execution of programs to prevent errors and improper use of the computer

Operating System Definition (Cont.)

- No universally accepted definition
- "Everything a vendor ships when you order an operating system" is good approximation

 But varies wildly
- "The one program running at all times on the computer" is the **kernel.** Everything else is either a system program (ships with the operating system) or an application program

Computer Startup

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

Computer System Organization

- Computer-system operation
 - One or more CPUs, device controllers connect through common bus providing access to shared



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 Timeline



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.

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.

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



Caching

- Important principle, performed at many levels in a computer (in hardware, operating system, software)
- Information in use copied from slower to faster storage temporarily
- Faster storage (cache) checked first to determine if information is there
 - If it is, information used directly from the cache (fast)
 - If not, data copied to cache and used there
- Cache smaller than storage being cached
 - Cache management important design problem
 - Cache size and replacement policy

Performance of Various Levels of Storage

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

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

Migration of Integer A from Disk to Register

 Multitasking environments must be careful to use most recent value, not matter where it is stored in the storage hierarchy



- Multiprocessor environment must provide cache coherency in hardware such that all CPUs have the most recent value in their cache
- Distributed environment situation even more complex
 - Several copies of a datum can exist
 - Various solutions covered in Chapter 17

Operating System Structure

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

Memory Layout for Multiprogrammed System



Operating-System Operations

- Interrupt driven by hardware
- Software error or request creates exception or trap

 Division by zero, request for operating system service
- Other process problems include infinite loop, processes modifying each other or the operating system
- **Dual-mode** operation allows OS to protect itself and other system components
 - User mode and kernel mode
 - Mode bit provided by hardware
 - Provides ability to distinguish when system is running user code or kernel code
 - Some instructions designated as privileged, only executable in kernel mode
 - System call changes mode to kernel, return from call resets it to user

Transition from User to Kernel Mode

- Timer to prevent infinite loop / process hogging resources
 - Set interrupt after specific period
 - Operating system decrements counter



Process Management

- A process is a program in execution. It is a unit of work within the system. Program is a *passive entity*, process is an *active entity*.
- Process needs resources to accomplish its task
 - CPU, memory, I/O, files
 - Initialization data
- Process termination requires reclaim of any reusable resources
- Single-threaded process has one **program counter** specifying location of next instruction to execute
 - Process executes instructions sequentially, one at a time, until completion
- Multi-threaded process has one program counter per thread
- Typically system has many processes, some user, some operating system running concurrently on one or more CPUs
 - Concurrency by multiplexing the CPUs among the processes / threads

Process Management Activities

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

- Creating and deleting both user and system processes
- Suspending and resuming processes
- Providing mechanisms for process synchronization
- Providing mechanisms for process communication
- Providing mechanisms for deadlock handling

Memory Management

- All data in memory before and after processing
- All instructions in memory in order to execute
- Memory management determines what is in memory when
 - Optimizing CPU utilization and computer response to users
- Memory management activities
 - Keeping track of which parts of memory are currently being used and by whom
 - Deciding which processes (or parts thereof) and data to move into and out of memory
 - Allocating and deallocating memory space as needed

Storage Management

- OS provides uniform, logical view of information storage
 - Abstracts physical properties to logical storage unit file
 - Each medium is controlled by device (i.e., disk drive, tape drive)
 - Varying properties include access speed, capacity, data-transfer rate, access method (sequential or random)
- File-System management
 - Files usually organized into directories
 - Access control on most systems to determine who can access what
 - OS activities include
 - Creating and deleting files and directories
 - Primitives to manipulate files and dirs
 - Mapping files onto secondary storage
 - Backup files onto stable (non-volatile) storage media

Mass-Storage Management

- Usually disks used to store data that does not fit in main memory or data that must be kept for a "long" period of time.
- Proper management is of central importance
- Entire speed of computer operation hinges on disk subsystem and its algorithms
- OS activities
 - Free-space management
 - Storage allocation
 - Disk scheduling
- Some storage need not be fast
 - Tertiary storage includes optical storage, magnetic tape
 - Still must be managed
 - Varies between WORM (write-once, read-many-times) and RW (read-write)

I/O Subsystem

- One purpose of OS is to hide peculiarities of hardware devices from the user
- I/O subsystem responsible for
 - Memory management of I/O including buffering (storing data temporarily while it is being transferred), caching (storing parts of data in faster storage for performance), spooling (the overlapping of output of one job with input of other jobs)
 - General device-driver interface
 - Drivers for specific hardware devices

Protection and Security

- Protection any mechanism for controlling access of processes or users to resources defined by the OS
- Security defense of the system against internal and external attacks
 - Huge range, including denial-of-service, worms, viruses, identity theft, theft of service
- Systems generally first distinguish among users, to determine who can do what
 - User identities (user IDs, security IDs) include name and associated number, one per user
 - User ID then associated with all files, processes of that user to determine access control
 - Group identifier (group ID) allows set of users to be defined and controls managed, then also associated with each process, file
 - Privilege escalation allows user to change to effective ID with more rights

Computing Environments

- Traditional computer
 - Blurring over time
 - Office environment
 - PCs connected to a network, terminals attached to mainframe or minicomputers providing batch and timesharing
 - Now portals allowing networked and remote systems access to same resources
 - Home networks
 - Used to be single system, then modems
 - Now firewalled, networked

Computing Environments (Cont.)

- Client-Server Computing
 - Dumb terminals supplanted by smart PCs
 - Many systems now **servers**, responding to requests generated by **clients**
 - Compute-server provides an interface to client to request services (i.e. database)
 - File-server provides interface for clients to store and retrieve files



Peer-to-Peer Computing

- Another model of distributed system
- P2P does not distinguish clients and servers
 - Instead all nodes are considered peers
 - May each act as client, server or both
 - Node must join P2P network
 - Registers its service with central lookup service on network, or
 - Broadcast request for service and respond to requests for service via *discovery protocol*
 - Examples include *Napster* and *Gnutella*

Web-Based Computing

- Web has become ubiquitous
- PCs most prevalent devices
- More devices becoming networked to allow web access
- New category of devices to manage web traffic among similar servers: load balancers
- Use of operating systems like Windows 95, client-side, have evolved into Linux and Windows XP, which can be clients and servers
Operating-System Structures

- Operating System Services
- User Operating System Interface
- System Calls
- Types of System Calls
- System Programs
- Operating System Design and Implementation
- Operating System Structure
- Virtual Machines
- Operating System Generation
- System Boot

Operating System Services

- One set of operating-system services provides functions that are helpful to the user:
 - User interface Almost all operating systems have a user interface (UI)
 - Varies between Command-Line (CLI), Graphics User Interface (GUI), Batch
 - Program execution The system must be able to load a program into memory and to run that program, end execution, either normally or abnormally (indicating error)
 - I/O operations A running program may require I/O, which may involve a file or an I/O device.
 - File-system manipulation The file system is of particular interest. Obviously, programs need to read and write files and directories, create and delete them, search them, list file Information, permission management.

Operating System Services (Cont.)

- One set of operating-system services provides functions that are helpful to the user (Cont):
 - Communications Processes may exchange information, on the same computer or between computers over a network
 - Communications may be via shared memory or through message passing (packets moved by the OS)
 - Error detection OS needs to be constantly aware of possible errors
 - May occur in the CPU and memory hardware, in I/O devices, in user program
 - For each type of error, OS should take the appropriate action to ensure correct and consistent computing
 - Debugging facilities can greatly enhance the user's and programmer's abilities to efficiently use the system

Operating System Services (Cont.)

- Another set of OS functions exists for ensuring the efficient operation of the system itself via resource sharing
 - Resource allocation When multiple users or multiple jobs running concurrently, resources must be allocated to each of them
 - Many types of resources Some (such as CPU cycles, mainmemory, and file storage) may have special allocation code, others (such as I/O devices) may have general request and release code.
 - Accounting To keep track of which users use how much and what kinds of computer resources
 - Protection and security The owners of information stored in a multiuser or networked computer system may want to control use of that information, concurrent processes should not interfere with each other
 - Protection involves ensuring that all access to system resources is controlled
 - **Security** of the system from outsiders requires user authentication, extends to defending external I/O devices from invalid access attempts
 - If a system is to be protected and secure, precautions must be instituted throughout it. A chain is only as strong as its weakest link.

User Operating System Interface - CLI

CLI allows direct command entry

- Sometimes implemented in kernel, sometimes by systems program
- Sometimes multiple flavors implemented shells
- Primarily fetches a command from user and executes it
 - Sometimes commands built-in, sometimes just names of programs
 - » If the latter, adding new features doesn't require shell modification

User Operating System Interface - GUI

- User-friendly **desktop** metaphor interface
 - Usually mouse, keyboard, and monitor
 - Icons represent files, programs, actions, etc
 - Various mouse buttons over objects in the interface cause various actions (provide information, options, execute function, open directory (known as a **folder**)
 - Invented at Xerox PARC
- Many systems now include both CLI and GUI interfaces
 - Microsoft Windows is GUI with CLI "command" shell
 - Apple Mac OS X as "Aqua" GUI interface with UNIX kernel underneath and shells available
 - Solaris is CLI with optional GUI interfaces (Java Desktop, KDE)

System Calls

- Programming interface to the services provided by the OS
- Typically written in a high-level language (C or C++)
- Mostly accessed by programs via a high-level Application Program Interface (API) rather than direct system call use
- Three most common APIs are Win32 API for Windows, POSIX API for POSIX-based systems (including virtually all versions of UNIX, Linux, and Mac OS X), and Java API for the Java virtual machine (JVM)
- Why use APIs rather than system calls?

(Note that the system-call names used throughout this text are generic)

Example of System Calls

 System call sequence to copy the contents of one <u>file to another file</u>

source file		destination file
	Example System Call Sequence Acquire input file name Write prompt to screen Accept input Acquire output file name Write prompt to screen Accept input Open the input file if file doesn't exist, abort Create output file if file exists, abort Loop Read from input file Write to output file Until read fails Close output file Write completion message to screen Terminate normally	

Example of Standard API

- Consider the ReadFile() function in the
- Win32 <u>API—a function for reading from a file</u>



- A description of the parameters passed to ReadFile()
 - HANDLE file—the file to be read
 - LPVOID buffer—a buffer where the data will be read into and written from
 - DWORD bytesToRead—the number of bytes to be read into the buffer
 - LPDWORD bytesRead—the number of bytes read during the last read
 - LPOVERLAPPED ovl—indicates if overlapped I/O is being used

System Call Implementation

- Typically, a number associated with each system call
 - System-call interface maintains a table indexed according to these numbers
- The system call interface invokes intended system call in OS kernel and returns status of the system call and any return values
- The caller need know nothing about how the system call is implemented
 - Just needs to obey API and understand what OS will do as a result call
 - Most details of OS interface hidden from programmer by API
 - Managed by run-time support library (set of functions built into libraries included with compiler)

API – System Call – OS Relationship



Standard C Library Example

• C program invoking printf() library call, which calls write() system call



System Call Parameter Passing

- Often, more information is required than simply identity of desired system call
 - Exact type and amount of information vary according to OS and call
- Three general methods used to pass parameters to the OS
 - Simplest: pass the parameters in *registers*
 - In some cases, may be more parameters than registers
 - Parameters stored in a *block*, or table, in memory, and address of block passed as a parameter in a register
 - This approach taken by Linux and Solaris
 - Parameters placed, or *pushed*, onto the *stack* by the program and *popped* off the stack by the operating system
 - Block and stack methods do not limit the number or length of parameters being passed

Parameter Passing via Table



Types of System Calls

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

MS-DOS execution



(a) At system startup (b) running a program

FreeBSD Running Multiple Programs

process D		
free memory		
process C		
interpreter		
process B		
kernel		

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

System Programs

- Provide a convenient environment for program development and execution
 - Some of them are simply user interfaces to system calls; others are considerably more complex
- File management Create, delete, copy, rename, print, dump, list, and generally manipulate files and directories
- Status information
 - Some ask the system for info date, time, amount of available memory, disk space, number of users
 - Others provide detailed performance, logging, and debugging information
 - Typically, these programs format and print the output to the terminal or other output devices
 - Some systems implement a registry used to store and retrieve configuration information

System Programs (cont'd)

- File modification
 - Text editors to create and modify files
 - Special commands to search contents of files or perform transformations of the text
- Programming-language support Compilers, assemblers, debuggers and interpreters sometimes provided
- Program loading and execution- Absolute loaders, relocatable loaders, linkage editors, and overlay-loaders, debugging systems for higher-level and machine language
- Communications Provide the mechanism for creating virtual connections among processes, users, and computer systems
 - Allow users to send messages to one another's screens, browse web pages, send electronic-mail messages, log in remotely, transfer files from one machine to another

Operating System Design and Implementation

- Design and Implementation of OS not "solvable", but some approaches have proven successful
- Internal structure of different Operating Systems can vary widely
- Start by defining goals and specifications
- Affected by choice of hardware, type of system
- User goals and System 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

Operating System Design and Implementation (Cont.)

- Important principle to separate
 Policy: What will be done?
 Mechanism: How to do it?
- 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

Simple 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



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

Layered Operating System



UNIX

- 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)				
Kernel	shells and commands compilers and interpreters system libraries				
	system-call interface to the kernel				
	signals terminal handling character I/O system terminal drivers	file system swapping block I/O system disk and tape drivers	CPU scheduling page replacement demand paging virtual memory		
	kernel interface to the hardware				
	terminal controllers terminals	device controllers disks and tapes	memory controllers physical memory		

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
- Detriments:
 - Performance overhead of user space to kernel space communication

Mac OS X Structure



Modules

- Most modern operating systems implement kernel modules
 - Uses object-oriented approach
 - Each core component is separate
 - Each talks to the others over known interfaces
 - Each is loadable as needed within the kernel
- Overall, similar to layers but with more flexible

Solaris Modular Approach



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



(a) Nonvirtual machine (b) virtual machine

Virtual Machines (Cont.)

- 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 operatingsystems 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
VMware Architecture

application	application	application	application
	guest operating system (free BSD) virtual CPU virtual memory virtual devices	guest operating system (Windows NT) virtual CPU virtual memory virtual devices virtualization layer	guest operating system (Windows XP) virtual CPU virtual memory virtual devices
\downarrow			
host operating system (Linux)			
hardware			
CPU memory I/O devices			

The Java Virtual Machine



Operating System Generation

- 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

System Boot

- Operating system must be made available to hardware so hardware can start it
 - Small piece of code **bootstrap loader**, locates the kernel, loads it into memory, and starts it
 - Sometimes two-step process where **boot block** at fixed location loads bootstrap loader
 - When power initialized on system, execution starts at a fixed memory location
 - Firmware used to hold initial boot code

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 in Memory



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)

process state

process number

program counter

registers

memory limits

list of open files

CPU Switch From Process to Process



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
- Processes migrate among 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) ⇒ (must be fast)
- Long-term scheduler is invoked very infrequently (seconds, minutes) ⇒ (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

Process Creation



C Program Forking Separate Process

```
int main()
Pid_t pid;
    /* fork another process */
    pid = fork();
    if (pid < 0) { /* error occurred */
            fprintf(stderr, "Fork Failed");
            exit(-1);
    else if (pid == 0) { /* child process */
            execlp("/bin/ls", "ls", NULL);
    else { /* parent process */
            /* parent will wait for the child to complete */
            wait (NULL);
            printf ("Child Complete");
            exit(0);
    }
}
```

A tree of processes on a typical Solaris



Process Termination

- Process executes last statement and asks the operating system to delete 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
 - If parent is exiting
 - Some operating system do not allow child to continue if its parent terminates
 - All children terminated *cascading termination*

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)



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 bidirectional

Module 2: Threads

- Overview
- Multithreading Models
- Threading Issues
- Pthreads
- Windows XP Threads
- Linux Threads
- Java Threads

Single and Multithreaded Processes



Benefits

- Responsiveness
- Resource Sharing
- Economy
- Utilization of MP Architectures

User Threads

- Thread management done by user-level threads library
- Three primary thread libraries:
 - POSIX Pthreads
 - Win32 threads
 - Java threads

Kernel Threads

- Supported by the Kernel
- Examples
 - Windows XP/2000
 - Solaris
 - Linux
 - Tru64 UNIX
 - Mac OS X

Multithreading Models

- Many-to-One
- One-to-One
- Many-to-Many

Many-to-One

- Many user-level threads mapped to single kernel thread
- Examples:
 - Solaris Green Threads
 - GNU Portable Threads

Many-to-One Model


One-to-One

- Each user-level thread maps to kernel thread
- Examples
 - Windows NT/XP/2000
 - Linux
 - Solaris 9 and later

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 prior to version 9
- Windows NT/2000 with the *ThreadFiber* package

Many-to-Many Model



Threading Issues

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

Semantics of fork() and exec()

 Does fork() duplicate only the calling thread or all threads?

Thread Cancellation

- Terminating a thread before it has finished
- Two general approaches:
 - Asynchronous cancellation terminates the target thread immediately
 - Deferred cancellation allows the target thread to periodically check if it should be cancelled

Signal Handling

- Signals are used in UNIX systems to notify a process that a particular event has occurred
- A signal handler is used to process signals
 - 1. Signal is generated by particular event
 - 2. Signal is delivered to a process
 - 3. Signal is handled
- Options:
 - Deliver the signal to the thread to which the signal applies
 - Deliver the signal to every thread in the process
 - Deliver the signal to certain threads in the process
 - Assign a specific thread to receive all signals for the process

Thread Pools

- Create a number of threads in a pool where they await work
- Advantages:
 - Usually slightly faster to service a request with an existing thread than create a new thread
 - Allows the number of threads in the application(s) to be bound to the size of the pool

Thread Specific Data

- Allows each thread to have its own copy of data
- Useful when you do not have control over the thread creation process (i.e., when using a thread pool)

Scheduler Activations

- Both M:M and Two-level models require communication to maintain the appropriate number of kernel threads allocated to the application
- Scheduler activations provide upcalls a communication mechanism from the kernel to the thread library
- This communication allows an application to maintain the correct number kernel threads

Thread libraries

Pthreads

Windows XP Threads

Linux Threads

Java Threads

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, Linux, Mac OS X)

Windows XP Threads

- Implements the one-to-one mapping
- Each thread contains
 - A thread id
 - Register set
 - Separate user and kernel stacks
 - Private data storage area
- The register set, stacks, and private storage area are known as the context of the threads
- The primary data structures of a thread include:
 - ETHREAD (executive thread block)
 - KTHREAD (kernel thread block)
 - TEB (thread environment block)

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 are managed by the JVM
- Java threads may be created by:
 - Extending Thread class
 - Implementing the Runnable interface

Java Thread States



CPU Scheduling

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Multiple-Processor Scheduling
- Real-Time Scheduling
- Thread Scheduling
- Operating Systems Examples
- Java Thread 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



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

<u>Process</u>	<u>Burst Time</u>
P_1	24
$\bar{P_2}$	3
P ₂	3

 Suppose that the processes arrive in the order: P₁, P₂, P₃ The Gantt Chart for the schedule is:

P ₁	P ₂	P ₃	
) 2	4 2	.7	30

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

FCFS Scheduling (Cont.)

Suppose that the processes arrive in the order

$$P_{2}, P_{3}, P_{1}$$

• The Gantt chart for the schedule is:



- 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	<u>Burst Time</u>
P_1	0.0	7
P_2	2.0	4
P_3	4.0	1
P_4	5.0	4

• SJF (non-preemptive)



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

Example of Preemptive SJF

Process	Arrival Time	<u>Burst Time</u>
P_1	0.0	7
P_2	2.0	4
P_3	4.0	1
P_4	5.0	4

• SJF (preemptive)



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

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
 - $q \text{ small} \Rightarrow q \text{ must}$ be large with respect to context switch, otherwise overhead is too high

Example of RR with Time Quantum = 20



• The Gantt chart is:



Typically, higher average turnaround than SJF, but better response

Time Quantum and Context Switch Time





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



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:
 - $-Q_0 RR$ with time quantum 8 milliseconds
 - $-Q_1 RR$ 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

Thread Scheduling

- Local Scheduling How the threads library decides which thread to put onto an available LWP
- Global Scheduling How the kernel decides which kernel thread to run next

Pthread Scheduling API

```
#include <pthread.h>
#include <stdio.h>
#define NUM THREADS 5
int main(int argc, char *argv[])
{
    int i:
    pthread t tid[NUM THREADS];
    pthread attr t attr;
    /* get the default attributes */
    pthread attr init(&attr);
    /* set the scheduling algorithm to PROCESS or SYSTEM */
    pthread attr setscope(&attr, PTHREAD SCOPE SYSTEM);
    /* set the scheduling policy - FIFO, RT, or OTHER */
    pthread attr setschedpolicy(&attr, SCHED OTHER);
    /* create the threads */
    for (i = 0; i < NUM THREADS; i++)
           pthread create(&tid[i],&attr,runner,NULL);
```

Pthread Scheduling API

```
/* now join on each thread */
for (i = 0; i < NUM THREADS; i++)
    pthread join(tid[i], NULL);
}
/* Each thread will begin control in this function */
void *runner(void *param)
{
    printf("I am a thread\n");
    pthread exit(0);
}</pre>
```

Operating System Examples

- Solaris scheduling
- Windows XP scheduling
- Linux scheduling

Solaris 2 Scheduling



Solaris Dispatch Table

priority	time quantum	time quantum expired	return from sleep		
0	200	0	50		
5	200	0	50		
10	160	0	51		
15	160	5	51		
20	120	10	52		
25	120	15	52		
30	80	20	53		
35	80	25	54		
40	40	30	55		
45	40	35	56		
50	40	40	58		
55	40	45	58		
59	20	49	59		

Windows XP 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

Linux Scheduling

- Two algorithms: time-sharing and real-time
- Time-sharing
 - Prioritized credit-based process with most credits is scheduled next
 - Credit subtracted when timer interrupt occurs
 - When credit = 0, another process chosen
 - When all processes have credit = 0, recrediting occurs
 - Based on factors including priority and history
- Real-time
 - Soft real-time
 - Posix.1b compliant two classes
 - FCFS and RR
 - Highest priority process always runs first

The Relationship Between Priorities and Time-slice length

numeric priority	relative priority		time quantum
0	highest		200 ms
•		real-time	
•		tasks	
•		lasks	
99			
100			
•		other	
•		tasks	
•		laono	
140	lowest		10 ms

List of Tasks Indexed According to Prorities



Algorithm Evaluation

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



5.08



In-5.7

	P _l	P ₂		P ₃	P ₄	P ₅	
0	1	0	39	94	2 4	9	61

In-5.8

	P ₃	P ₄	P ₁	P ₅	P ₂
(0 3	3 1	.0 2	0 3	2 61

In-5.9

	P ₁	P ₂	P ₃	\mathbb{P}_4	P ₅	P ₂	P ₅	P ₂	
С) 1	0	20 2	3 3	0 4	0	50 52	2	61

Dispatch Latency



Process Synchronization

- Background
- The Critical-Section Problem
- Peterson's Solution
- Synchronization Hardware
- Semaphores
- Classic Problems of Synchronization
- Monitors
- Synchronization Examples
- Atomic Transactions

Background

- Concurrent access to shared data may result in data inconsistency
- Maintaining data consistency requires mechanisms to ensure the orderly execution of cooperating processes is called as race condition.
- Suppose that we wanted to provide a solution to the consumer-producer problem that fills all the buffers. We can do so by having an integer count that keeps track of the number of full buffers. Initially, count is set to 0. It is incremented by the producer after it produces a new buffer and is decremented by the consumer after it consumes a buffer.

Solution to Critical-Section Problem

- 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

Peterson's Solution

- Two process solution
- Assume that the LOAD and STORE instructions are atomic; that is, cannot be interrupted.
- The two processes share two variables:
 - int turn;
 - Boolean flag[2]
- The variable turn indicates whose turn it is to enter the critical section.
- The flag array is used to indicate if a process is ready to enter the critical section. flag[i] = true implies that process P_i is ready!

Algorithm for Process P_i

```
do {
```

```
flag[i] = TRUE;
turn = j;
while ( flag[j] && turn == j);
```

```
CRITICAL SECTION
```

flag[i] = FALSE;

REMAINDER SECTION

} while (TRUE);

Synchronization Hardware

- Many systems provide hardware support for critical section code
- Uniprocessors could disable interrupts
 - Currently running code would execute without preemption
 - Generally too inefficient on multiprocessor systems
 - Operating systems using this not broadly scalable
- Modern machines provide special atomic hardware instructions
 - Atomic = non-interruptable
 - Either test memory word and set value
 - Or swap contents of two memory words

TestAndndSet Instruction

• Definition:

```
boolean TestAndSet (boolean
*target)
{
    boolean rv = *target;
    *target = TRUE;
    return rv:
}
```

Solution using TestAndSet

- Shared boolean variable lock., initialized to false.
- Solution:
 - do {
 while (TestAndSet (&lock))
 ; /* do nothing
 - // critical section

lock = FALSE;

// remainder section

} while (TRUE);

Swap Instruction

• Definition:

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

Solution using Swap

- Shared Boolean variable lock initialized to FALSE; Each process has a local Boolean variable key.
- Solution:

```
do {
	key = TRUE;
	while ( key == TRUE)
	Swap (&lock, &key );
```

// critical section

lock = FALSE;

// remainder section

} while (TRUE);

Semaphore

- Synchronization tool that does not require busy waiting
- Semaphore *S* integer variable
- Two standard operations modify S: wait() and signal()

Originally called P() and V()

- Less complicated
- Can only be accessed via two indivisible (atomic) operations

```
- wait (S) {
    while S <= 0
    ; // no-op
    S--;
    }
- signal (S) {
    S++;
    }
</pre>
```

Semaphore as General Synchronization Tool

- 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

Also known as mutex locks

- Can implement a counting semaphore S as a binary semaphore
- Provides mutual exclusion
 - Semaphore S; // initialized to 1
 - wait (S);

Critical Section

signal (S);

Semaphore Implementation

- Must guarantee that no two processes can execute wait

 () and signal () on the same semaphore at the same time
- Thus, implementation becomes the critical section problem where the wait and signal code are placed in the crtical section.
 - Could now have busy waiting in critical section implementation
 - But implementation code is short
 - Little busy waiting if critical section rarely occupied
- Note that applications may spend lots of time in critical sections and therefore this is not a good solution.
Semaphore Implementation with no Busy waiting

- With each semaphore there is an associated waiting queue. Each entry in a waiting queue has two data items:
 - value (of type integer)
 - pointer to next record in the list
- Two operations:
 - block place the process invoking the operation on the appropriate waiting queue.
 - wakeup remove one of processes in the waiting queue and place it in the ready queue.

Semaphore Implementation with no Busy waiting (Cont.)

• Implementation of wait:

```
wait (S){
    value--;
    if (value < 0) {
        add this process to waiting queue
        block(); }
}</pre>
```

• Implementation of signal:

```
Signal (S){
    value++;
    if (value <= 0) {
        remove a process P from the waiting queue
        wakeup(P); }
}</pre>
```

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 o be two semaphores initialized to 1



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

Classical Problems of Synchronization

- Bounded-Buffer Problem
- Readers and Writers Problem
- Dining-Philosophers Problem

Bounded-Buffer Problem

- *N* buffers, each can hold one item
- Semaphore mutex initialized to the value 1
- Semaphore full initialized to the value 0
- Semaphore empty initialized to the value N.

Bounded Buffer Problem (Cont.)

• The structure of the producer process

```
do {
    // produce an item
    wait (empty);
```

```
wait (mutex);
```

```
// add the item to the buffer
```

```
signal (mutex);
signal (full);
} while (true);
```

Bounded Buffer Problem (Cont.)

• The structure of the consumer process

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

// remove an item from buffer

signal (mutex);
signal (empty);

// consume the removed item

} while (true);

Readers-Writers Problem

- A data set is shared among a number of concurrent processes
 - Readers only read the data set; they do not perform any updates
 - Writers can both read and write.
- Problem allow multiple readers to read at the same time. Only one single writer can access the shared data at the same time.
- Shared Data
 - Data set
 - Semaphore mutex initialized to 1.
 - Semaphore wrt initialized to 1.
 - Integer readcount initialized to 0.

Readers-Writers Problem (Cont.)

• The structure of a writer process

```
do {
    wait (wrt) ;
    // writing is performed
    signal (wrt) ;
} while (true)
```

Readers-Writers Problem (Cont.)

• The structure of a reader process

```
do {
    wait (mutex);
    readcount ++;
    if (readercount == 1) wait (wrt);
    signal (mutex)
```

// reading is performed

```
wait (mutex) ;
readcount --;
if redacount == 0) signal (wrt) ;
signal (mutex) ;
} while (true)
```

Dining-Philosophers Problem



- Shared data
 - Bowl of rice (data set)
 - Semaphore chopstick [5] initialized to 1

Dining-Philosophers Problem (Cont.)

• The structure of Philosopher *i*:

```
Do {
   wait ( chopstick[i] );
   wait ( chopStick[ (i + 1) % 5] );
       // eat
   signal ( chopstick[i] );
   signal (chopstick[ (i + 1) % 5] );
      // think
} while (true) ;
```

Problems with Semaphores

- Incorrect use of semaphore operations:
 - signal (mutex) wait (mutex)
 - wait (mutex) … wait (mutex)
 - Omitting of wait (mutex) or signal (mutex) (or both)

Monitors

- A high-level abstraction data type that provides a convenient and effective mechanism for process synchronization
- Only one process may be active within the monitor at a time

```
monitor monitor-name
{
    // shared variable declarations
    procedure P1 (...) { .... }
    ...
    procedure Pn (...) { .....}
    Initialization code ( ....) { .... }
    ...
    }
}
```

Schematic view of a Monitor



Condition Variables

- condition x, y;
- Two operations on a condition variable:
 - x.wait () a process that invokes the operation is

suspended.

x.signal () – resumes one of processes (if any)
 tha

invoked x.wait ()

Monitor with Condition Variables



Solution to Dining Philosophers

```
monitor DP
 {
    enum { THINKING; HUNGRY, EATING) state [5] ;
    condition self [5];
    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);
    }
```

Solution to Dining Philosophers (cont)

```
void test (int i) {
    if ( (state[(i + 4) % 5] != EATING) &&
        (state[i] == HUNGRY) &&
        (state[(i + 1) % 5] != EATING) ) {
            state[i] = EATING ;
            self[i].signal () ;
        }
    }
initialization_code() {
    for (int i = 0; i < 5; i++)
        state[i] = THINKING;
}</pre>
```

}

Module 3: 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.
 - P₁ and P₂ each hold one tape drive and each needs another one.
- Example
 - semaphores A and B, initialized to 1

 P0
 P1

 wait (A);
 wait(B)

 wait (B);
 wait(A)

Operating System Concepts

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 upif 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 {P₀, P₁, ..., P₀} of waiting processes such that P₀ is waiting for a resource that is held by P₁, P₁ is waiting for a resource that is held by P₂, ..., P_{n-1} is waiting for a resource that is held by P₀.

Resource-Allocation Graph

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

- V is partitioned into two types:
 - $P = \{P_1, P_2, ..., P_n\}$, the set consisting of all the processes in the system.
 - $R = \{R_1, R_2, ..., R_m\}$, the set consisting of all resource types in the system.
- request edge directed edge $P_1 \rightarrow R_i$
- 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_i

• P_i is holding an instance of R_i







Operating System Concepts

Resource Allocation Graph With A Deadlock R_3 R_1 P_3 P. P_2 R_2 R_4

Operating System Concepts

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 <P₁, P₂, ..., 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_i, with j<I.
 - If P_i resource needs are not immediately available, then P_i can wait until all P_i 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.

Safe, unsafe, deadlock state spaces



Operating String System of the cepts

Safe State

To illustrate, we consider a system with twelve magnetic tape drives and three processes: P0, P1, and P2. Process P0 requires ten tape drives, process P1 may need as many as four tape drives, and process P2 may need up to nine tape drives. Suppose that, at time t0, process P0 is holding five tape drives, process P1 is holding two tape drives, and process P2 is holding two tape drives. (Thus, there are three free tape drives.)

Maximum Needs		Current Needs	maximum need to finished
Po	10	5	5
P_1	4	2	2
P_2	9	2	7

Safe State

•At time t0, the system is in a safe state. The sequence **<P1,P0,P2>** satisfies the safety condition. Process P1 can immediately be allocated all its tape drives and then return them (the system will then have five available tape drives);

•then process PO can get all its tape drives and return them (the system will then have ten available tape drives);

•and finally process P2 can get all its tape drives and return them (the system will then have all twelve tape drives available).

•A system can go from a safe state to an unsafe state. Suppose that, at time t1, process P2 requests and is allocated one more tape drive.

•The system is no longer in a safe state. At this point, only process P1 can be allocated all its tape drives. When it returns them, the system will have only four available tape drives.

• Since process P0 is allocated five tapes drives but has a maximum of ten, it may request five more tape drives.

• If it does so, it will have to wait, because they are unavailable. Similarly, process P2 may request six additional tape drives and have to wait, resulting in a deadlock.

•Our mistake was in granting the request from process P2 for one more tape drive.

•If we had made P2 wait until either of the other processes had finished and released its resources, then we could have avoided the deadlock

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.
Resource-Allocation Graph For Deadlock Avoidance



Unsafe State In A Resource-Allocation Graph R P_2

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_i 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_i.
- Allocation: n x m matrix. If Allocation[i,j] = k then P_i is currently allocated k instances of R_i.
- Need: n x m matrix. If Need[i,j] = k, then P_i may need k more instances of R_i 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.

- 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_i = request vector for process P_i . If Request_i[j] = k then process P_i wants k instances of resource type R_i .

- 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 P₀ through P₄; 3 resource types A (10 instances), B (5instances, and C (7 instances).
- Snapshot at time T_0 :

	<u>Allocation</u>	<u>Max</u>	<u>Available</u>
	A B C	A B C	A B C
P ₀	010	753	332
P_1	200	322	
<i>P</i> ₂	302	902	
<i>P</i> ₃	211	222	
P_4	002	433	

Example (Cont.)

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

	<u>Need</u>	
	ABC	
<i>P</i> ₀	743	
<i>P</i> ₁	122	
<i>P</i> ₂	600	
<i>P</i> ₃	011	
<i>P</i> ₄	431	

The system is in a safe state since the sequence < P₁, P₃, P₄, P₂, P₀> satisfies safety criteria.

Example (Cont.): P₁ request (1,0,2)

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

	<u>Allocation</u>	<u>Need</u>	<u>Available</u>
	ABC	ABC	A B C
P ₀	010	743	230
P_1	302	020	
P ₂	301	600	
P ₃	211	011	
P_4	002	431	

- Executing safety algorithm shows that sequence <P₁, P₃, P₄, P₀, P₂> 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 acycle 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



Corresponding wait-for graph

 P_5

P,

P

(b)

3

Operating System Concepts

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* [*ij*] = *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) $Request_i \leq Work$

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*, $1 \le i \le n$, then the system is in deadlock state. Moreover, if *Finish*[*i*] = *false*, then *P*_{*i*} is deadlocked.

Algorithm requires an order of $m \times n^2$ operations to detect whether the system is in deadlocked state.

Example of Detection Algorithm

- Five processes P₀ through P₄; three resource types
 A (7 instances), B (2 instances), and C (6 instances).
- Snapshot at time T_0 :

	<u>Allocation</u>	Request	<u>Available</u>
	A B C	A B C	A B C
P ₀	010	000	000
P_1	200	202	
P ₂	303	000	
P ₃	211	100	
P_4	002	002	

• Sequence $\langle P_0, P_2, P_3, P_1, P_4 \rangle$ will result in *Finish*[*i*] = true for all *i*.

Example (Cont.)

• *P*₂ requests an additional instance of type *C*.

	<u>Request</u>	
	ABC	
P ₀	000	
<i>P</i> ₁	201	
2 2	001	
7 3	100	
D ₄	002	

- State of system?
 - Can reclaim resources held by process P₀, 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 fro 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.

Memory Management

- Background
- Logical versus Physical Address Space
- Swapping
- Contiguous Allocation
- Paging
- Segmentation
- Segmentation with Paging

Background

- Program must be brought into memory and placed within a process for it to be executed.
- *Input queue* collection of processes on the disk that are waiting to be brought into memory for execution.
- User programs go through several steps before being executed.

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*).

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.

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

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.

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 and Microsoft Windows.

Schematic View of Swapping



Operating System Concepts

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.

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)



Operating System Concepts

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 entier 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



Operating System Concepts

Paging Example



Operating System Concepts

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 registers* or *translation look-aside buffers (TLBs)*

Associative Register

• Associative registers – parallel search

	Page #	Frame #
Address translation (A ['] ,		

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

Effective Access Time

- Associative Lookup = ε time unit
- Assume memory cycle time is 1 microsecond
- Hit ration percentage of times that a page number is found in the associative registers; ration related to number of associative registers.
- Hit ratio = α
- 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.

Two-Level Page-Table Scheme



Operating System Concepts

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

 where p_i is an index into the outer page table, and p₂ is the displacement within the page of the outer page table.

Address-Translation Scheme

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



Multilevel Paging and Performance

- Since each level is stored as a separate table in memory, covering a logical address to a physical one may take four memory accesses.
- Even though time needed for one memory access is quintupled, caching permits performance to remain reasonable.
- Cache hit rate of 98 percent yields:

effective access time = 0.98 x 120 + 0.02 x 520

= 128 nanoseconds.

which is only a 28 percent slowdown in memory access time.

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



Operating System Concepts

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,

local variables, global variables,

common block,

stack,

symbol table, arrays

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

Sharing of segments



Operating System Concepts

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



Module 4: Virtual Memory

- Background
- Demand Paging
- Performance of Demand Paging
- Page Replacement
- Page-Replacement Algorithms
- Allocation of Frames
- Thrashing
- Other Considerations
- Demand Segmenation

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.
 - Need to allow pages to be swapped in and out.
- Virtual memory can be implemented via:
 - Demand paging
 - Demand segmentation

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

Valid-Invalid Bit

- With each page table entry a valid–invalid bit is associated (1 ⇒ in-memory, 0 ⇒ not-in-memory)
- Initially valid—invalid but is set to 0 on all entries.
- Example of a page table snapshot.



Page Fault

- If there is ever a reference to a page, first reference will trap to OS ⇒ 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



Operating System Concepts

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 $0 \le p \le 1.0$
 - if p = 0 no page faults
 - if p = 1, every reference is a fault
- Effective Access Time (EAT)
 - EAT = (1 p) x memory access
 - + p (page fault overhead
 - + [swap page out]
 - + swap page in
 - + 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 (in msec)

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.

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.

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)



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



- How do you know this?
- Used for measuring how well your algorithm performs.

Least Recently Used (LRU) Algorithm

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



- 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 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

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.

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 $p_i = \frac{s_i}{S} \times m$
 $m = 64$
 $s_i = 10$
 $s_2 = 127$
 $a_1 = \frac{10}{137} \times 64 \approx 5$
 $a_2 = \frac{127}{137} \times 64 \approx 59$

Operating System Concepts

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 Diagram



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

File-System Interface

- File Concept
- Access :Methods
- Directory Structure
- Protection
- Consistency Semantics

File Concept

- Contiguous logical address space
- Types:
 - Data
 - numeric
 - character
 - binary
 - Program

File Structure

- None sequence of words, bytes
- Simple record structure
 - Lines
 - Fixed length
 - Variable length
- Complex Structures
 - Formatted document
 - Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters.
- Who decides:
 - Operating system
 - Program

File Attributes

- **Name** only information kept in human-readable form.
- **Type** needed for systems that support different types.
- **Location** pointer to file location on device.
- **Size** current file size.
- **Protection** controls who can do reading, writing, executing.
- **Time**, **date**, **and user identification** data for protection, security, and usage monitoring.
- Information about files are kept in the directory structure, which is maintained on the disk.

File Operations

- create
- write
- read
- reposition within file file seek
- delete
- truncate
- open(F_i) search the directory structure on disk for entry F_i , and move the content of entry to memory.
- close (F_i) move the content of entry F_i in memory to directory structure on disk.

File Types – name, extension

File Type	Usual extension	Function	
Executable	exe, com, bin or none	ready-to-run machine- language program	
Object	obj, o	complied, machine language, not linked	
Source code	c, p, pas, 177, asm, a	source code in various languages	
Batch	bat, sh	commands to the command interpreter	
Text	txt, doc	textual data documents	
Word processor	wp, tex, rrf, etc.	various word-processor formats	
Library	lib, a	libraries of routines	
Print or view	ps, dvi, gif	ASCII or binary file	
Archive	arc, zip, tar	related files grouped into one file, sometimes compressed.	

Access Methods

• Sequential Access

read next write next reset no read after last write (rewrite)

Direct Access

read n write n position to n read next write next rewrite n

n = relative block number

Operating System Concepts

Directory Structure

• A collection of nodes containing information about all files.



- Both the directory structure and the files reside on disk.
- Backups of these two structures are kept on tapes.

Operating System Concepts

Information in a Device Directory

- Name
- Type
- Address
- Current length
- Maximum length
- Date last accessed (for archival)
- Date last updated (for dump)
- Owner ID (who pays)
- Protection information (discuss later)

Operations Performed on Director

- Search for a file
- Create a file
- Delete a file
- List a directory
- Rename a file
- Traverse the file system

Organize the Directory (Logically) to Obtain

- Efficiency locating a file quickly.
- Naming convenient to users.
 - Two users can have same name for different files.
 - The same file can have several different names.
- Grouping logical grouping of files by properties, (e.g., all Pascal programs, all games, ...)

Single-Level Directory

• A single directory for all users.



- Naming problem
- Grouping problem

Two-Level Directory

• Separate directory for each user.



- Path name
- Can have the saem file name for different user
- Efficient searching
- No grouping capability

Operating System Concepts

Tree-Structured Directories



Operating System Concepts

Tree-Structured Directories (Cont.)

- Efficient searching
- Grouping Capability
- Current directory (working directory)
 - cd /spell/mail/prog
 - type list

Tree-Structured Directories (Cont.)

- Absolute or relative path name
- Creating a new file is done in current directory.
- Delete a file

rm <file-name>

• Creating a new subdirectory is done in current directory.

mkdir <dir-name>

Example: if in current directory /spell/mail

mkdir count



• Deleting "mail" \Rightarrow deleting the entire subtree rooted by "mail".

Operating System Concepts

Acyclic-Graph Directories

• Have shared subdirectories and files.



Acyclic-Graph Directories (Cont.)

- Two different names (aliasing)
- If *dict* deletes *list* ⇒ dangling pointer.
 Solutions:
 - Backpointers, so we can delete all pointers.
 Variable size records a problem.
 - Backpointers using a daisy chain organization.
 - Entry-hold-count solution.

General Graph Directory



General Graph Directory (Cont.)

- How do we guarantee no cycles?
 - Allow only links to file not subdirectories.
 - Garbage collection.
 - Every time a new link is added use a cycle detection algorithm to determine whether it is OK.

Protection

- File owner/creator should be able to control:
 - what can be done
 - by whom
- Types of access
 - Read
 - Write
 - Execute
 - Append
 - Delete
 - List

Access Lists and Groups

- Mode of access: read, write, execute
- Three classes of users

			RVVX
a) owner access	7	\Rightarrow	111
			RWX
b) groups access	6	\Rightarrow	110
			RWX
c) public access	1	\Rightarrow	001

- Ask manager to create a group (unique name), say G, and add some users to the group.
- For a particular file (say *game*) or subdirectory, define an appropriate access.



• Attach a group to a file

chgrp G game

Operating System Concepts

File-System Implementation

- File-System Structure
- Allocation Methods
- Free-Space Management
- Directory Implementation
- Efficiency and Performance
- Recovery

File-System Structure

- File structure
 - Logical storage unit
 - Collection of related information
- File system resides on secondary storage (disks).
- File system organized into layers.
- *File control block* storage structure consisting of information about a file.

Contiguous Allocation

- Each file occupies a set of contiguous blocks on the disk.
- Simple only starting location (block #) and length (number of blocks) are required.
- Random access.
- Wasteful of space (dynamic storage-allocation problem).
- Files cannot grow.
- Mapping from logical to physical.



- Block to be accessed = ! + starting address
- Displacement into block = R

Linked Allocation

• Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk.



• Allocate as needed, link together; e.g., file starts at block 9


Linked Allocation (Cont.)

- Simple need only starting address
- Free-space management system no waste of space
- No random access
- Mapping



- Block to be accessed is the Qth block in the linked chain of blocks representing the file.
- Displacement into block = R + 1
- File-allocation table (FAT) disk-space allocation used by MS-DOS and OS/2.

Indexed Allocation

- Brings all pointers together into the *index block*.
- Logical view.



index table

Example of Indexed Allocation



Indexed Allocation (Cont.)

- Need index table
- Random access
- Dynamic access without external fragmentation, but have overhead of index block.
- Mapping from logical to physical in a file of maximum size of 256K words and block size of 512 words. We need only 1 block for index table.

- Q = displacement into index table
- R = displacement into block

Indexed Allocation – Mapping (Cont.)

- Mapping from logical to physical in a file of unbounded length (block size of 512 words).
- Linked scheme Link blocks of index table (no limit on size).



- Q_2 = displacement into block of index table
- R_2 displacement into block of file:

Indexed Allocation – Mapping (Cont.)

• Two-level index (maximum file size is 512³)



- Q_1 = displacement into outer-index



- Q_2 = displacement into block of index table
- R_2 displacement into block of file:

Indexed Allocation – Mapping (Cont.)



Combined Scheme: UNIX (4K bytes per block)



Free-Space Management

• Bit vector (*n* blocks)



bit[*i*] =
$$\begin{cases} 0 \Rightarrow block[i] free \\ 1 \Rightarrow block[i] occupied \end{cases}$$

Block number calculation

(number of bits per word) * (number of 0-value words) + offset of first 1 bit

Free-Space Management (Cont.)

• Bit map requires extra space. Example:

block size = 2^{12} bytes disk size = 2^{30} bytes (1 gigabyte) $n = 2^{30}/2^{12} = 2^{18}$ bits (or 32K bytes)

- Easy to get contiguous files
- Linked list (free list)
 - Cannot get contiguous space easily
 - No waste of space
- Grouping
- Counting

Free-Space Management (Cont.)

- Need to protect:
 - Pointer to free list
 - Bit map
 - Must be kept on disk
 - Copy in memory and disk may differ.
 - Cannot allow for block[i] to have a situation where bit[i] = 1 in memory and bit[i] = 0 on disk.
 - Solution:
 - Set bit[*i*] = 1 in disk.
 - Allocate block[*i*]
 - Set bit[*i*] = 1 in memory

Directory Implementation

- Linear list of file names with pointer to the data blocks.
 - simple to program
 - time-consuming to execute
- Hash Table linear list with hash data structure.
 - decreases directory search time
 - collisions situations where two file names hash to the same location
 - fixed size

Efficiency and Performance

- Efficiency dependent on:
 - disk allocation and directory algorithms
 - types of data kept in file's directory entry
- Performance
 - disk cache separate section of main memory for frequently sued blocks
 - free-behind and read-ahead techniques to optimize sequential access
 - improve PC performance by dedicating section of memroy as virtual disk, or RAM disk.

Various Disk-Caching Locations



Recovery

- Consistency checker compares data in directory structure with data blocks on disk, and tries to fix inconsistencies.
- Use system programs to *back up* data from disk to another storage device (floppy disk, magnetic tape).
- Recover lost file or disk by *restoring* data from backup.

Module 5: Mass-Storage Systems

- Overview of Mass Storage Structure
- Disk Structure
- Disk Attachment
- Disk Scheduling
- Disk Management
- Swap-Space Management
- RAID Structure
- Stable-Storage Implementation

Overview of Mass Storage Structure

- Magnetic disks provide bulk of secondary storage of modern computers
 - Drives rotate at 60 to 250 times per second
 - Transfer rate is rate at which data flow between drive and computer
 - Positioning time (random-access time) is time to move disk arm to desired cylinder (seek time) and time for desired sector to rotate under the disk head (rotational latency)
 - Head crash results from disk head making contact with the disk surface --That's bad
- Disks can be removable
- Drive attached to computer via I/O bus
 - Busses vary, including EIDE, ATA, SATA, USB, Fibre Channel, SCSI, SAS, Firewire
 - Host controller in computer uses bus to talk to disk controller built into drive or storage array

Moving-head Disk Mechanism



Hard Disks

- Platters range from .85" to 14" (historically)
 - Commonly 3.5", 2.5", and 1.8"
- Range from 30GB to 3TB per drive
- Performance
 - Transfer Rate theoretical 6 Gb/sec
 - Effective Transfer Rate real 1Gb/sec
 - Seek time from 3ms to 12ms 9ms common for desktop drives
 - Average seek time measured or calculated based on 1/3 of tracks
 - Latency based on spindle speed
 - 1 / (RPM / 60) = 60 / RPM
 - Average latency = ½ latency

Spindle [rpm]	Average latency [ms]
4200	7.14
5400	5.56
7200	4.17
10000	3
15000	2

(From Wikipedia)

Hard Disk Performance

- Access Latency = Average access time = average seek time + average latency
 - For fastest disk 3ms + 2ms = 5ms
 - For slow disk 9ms + 5.56ms = 14.56ms
- Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead
- For example to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a .1ms controller overhead =
 - 5ms + 4.17ms + 0.1ms + transfer time =
 - Transfer time = 4KB / 1Gb/s * 8Gb / GB * 1GB / 1024²KB = 32 / (1024²) = 0.031 ms
 - Average I/O time for 4KB block = 9.27ms + .031ms = 9.301ms

Solid-State Disks

- Nonvolatile memory used like a hard drive
 - Many technology variations
- Can be more reliable than HDDs
- More expensive per MB
- Maybe have shorter life span
- Less capacity
- But much faster
- Busses can be too slow -> connect directly to PCI for example
- No moving parts, so no seek time or rotational latency

Magnetic Tape

- Was early secondary-storage medium
 - Evolved from open spools to cartridges
- Relatively permanent and holds large quantities of data
- Access time slow
- Random access ~1000 times slower than disk
- Mainly used for backup, storage of infrequently-used data, transfer medium between systems
- Kept in spool and wound or rewound past read-write head
- Once data under head, transfer rates comparable to disk
 - 140MB/sec and greater
- 200GB to 1.5TB typical storage
- Common technologies are LTO-{3,4,5} and T10000

Disk Structure

- Disk drives are addressed as large 1-dimensional arrays of logical blocks, where the logical block is the smallest unit of transfer
 - Low-level formatting creates logical blocks on physical media
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - Sector 0 is the first sector of the first track on the outermost cylinder
 - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost
 - Logical to physical address should be easy
 - Except for bad sectors
 - Non-constant # of sectors per track via constant angular velocity

Disk Attachment

- Host-attached storage accessed through I/O ports talking to I/O busses
- SCSI itself is a bus, up to 16 devices on one cable, SCSI initiator requests operation and SCSI targets perform tasks
 - Each target can have up to 8 logical units (disks attached to device controller)
- FC is high-speed serial architecture
 - Can be switched fabric with 24-bit address space the basis of storage area networks (SANs) in which many hosts attach to many storage units
- I/O directed to bus ID, device ID, logical unit (LUN)

Storage Array

- Can just attach disks, or arrays of disks
- Storage Array has controller(s), provides features to attached host(s)
 - Ports to connect hosts to array
 - Memory, controlling software (sometimes NVRAM, etc)
 - A few to thousands of disks
 - RAID, hot spares, hot swap (discussed later)
 - Shared storage -> more efficiency
 - Features found in some file systems
 - Snaphots, clones, thin provisioning, replication, deduplication, etc

Storage Area Network

- Common in large storage environments
- Multiple hosts attached to multiple storage arrays flexible



Storage Area Network (Cont.)

- SAN is one or more storage arrays
 - Connected to one or more Fibre Channel switches
- Hosts also attach to the switches
- Storage made available via LUN Masking from specific arrays to specific servers
- Easy to add or remove storage, add new host and allocate it storage
 - Over low-latency Fibre Channel fabric
- Why have separate storage networks and communications networks?
 - Consider iSCSI, FCOE

Network-Attached Storage

- Network-attached storage (NAS) is storage made available over a network rather than over a local connection (such as a bus)
 - Remotely attaching to file systems
- NFS and CIFS are common protocols
- Implemented via remote procedure calls (RPCs) between host and storage over typically TCP or UDP on IP network
- **iSCSI** protocol uses IP network to carry the SCSI protocol
 - Remotely attaching to devices (blocks)



Disk Scheduling

- The operating system is responsible for using hardware efficiently for the disk drives, this means having a fast access time and disk bandwidth
- Minimize seek time
- Seek time \approx seek distance
- Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer

Disk Scheduling (Cont.)

- There are many sources of disk I/O request
 - OS
 - System processes
 - Users processes
- I/O request includes input or output mode, disk address, memory address, number of sectors to transfer
- OS maintains queue of requests, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
 - Optimization algorithms only make sense when a queue exists

Disk Scheduling (Cont.)

- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying "depth")
- Several algorithms exist to schedule the servicing of disk I/O requests
- The analysis is true for one or many platters
- We illustrate scheduling algorithms with a request queue (0-199)

98, 183, 37, 122, 14, 124, 65, 67

Head pointer 53

FCFS

Illustration shows total head movement of 640 cylinders



SSTF

- Shortest Seek Time First selects the request with the minimum seek time from the current head position
- SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests
- Illustration shows total head movement of 236 cylinders



SCAN

- The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues.
- SCAN algorithm Sometimes called the elevator algorithm
- Illustration shows total head movement of 236 cylinders
- But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest

SCAN (Cont.)



C-SCAN

- Provides a more uniform wait time than SCAN
- The head moves from one end of the disk to the other, servicing requests as it goes
 - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
- Treats the cylinders as a circular list that wraps around from the last cylinder to the first one
- Total number of cylinders?
C-SCAN (Cont.)



C-LOOK

- LOOK a version of SCAN, C-LOOK a version of C-SCAN
- Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk
- Total number of cylinders?

C-LOOK (Cont.)



Selecting a Disk-Scheduling Algorithm

- SSTF is common and has a natural appeal
- SCAN and C-SCAN perform better for systems that place a heavy load on the disk
 - Less starvation
- Performance depends on the number and types of requests
- Requests for disk service can be influenced by the file-allocation method
 - And metadata layout
- The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary
- Either SSTF or LOOK is a reasonable choice for the default algorithm
- What about rotational latency?
 - Difficult for OS to calculate
- How does disk-based queueing effect OS queue ordering efforts?

Disk Management

- Low-level formatting, or physical formatting Dividing a disk into sectors that the disk controller can read and write
 - Each sector can hold header information, plus data, plus error correction code (ECC)
 - Usually 512 bytes of data but can be selectable
- To use a disk to hold files, the operating system still needs to record its own data structures on the disk
 - Partition the disk into one or more groups of cylinders, each treated as a logical disk
 - Logical formatting or "making a file system"
 - To increase efficiency most file systems group blocks into clusters
 - Disk I/O done in blocks
 - File I/O done in clusters

Disk Management (Cont.)

- Raw disk access for apps that want to do their own block management, keep OS out of the way (databases for example)
- Boot block initializes system
 - The bootstrap is stored in ROM
 - **Bootstrap loader** program stored in boot blocks of boot partition
- Methods such as **sector sparing** used to handle bad blocks

Booting from a Disk in Windows



Swap-Space Management

- Swap-space Virtual memory uses disk space as an extension of main memory
 - Less common now due to memory capacity increases
- Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition (raw)
- Swap-space management
 - 4.3BSD allocates swap space when process starts; holds text segment (the program) and data segment
 - Kernel uses swap maps to track swap-space use
 - Solaris 2 allocates swap space only when a dirty page is forced out of physical memory, not when the virtual memory page is first created
 - File data written to swap space until write to file system requested
 - Other dirty pages go to swap space due to no other home
 - Text segment pages thrown out and reread from the file system as needed
- What if a system runs out of swap space?
- Some systems allow multiple swap spaces