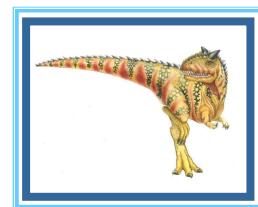


Chapter 1: Introduction



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

1



Chapter 1: Introduction

- What Operating Systems Do
- Computer-System Organization
- Computer-System Architecture
- Operating-System Operations
- Resource Management
- Security and Protection
- Virtualization



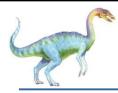
Operating System Concepts – 10th Edition

1.2

Silberschatz, Galvin and Gagne ©2018

2

1



What is an Operating System?

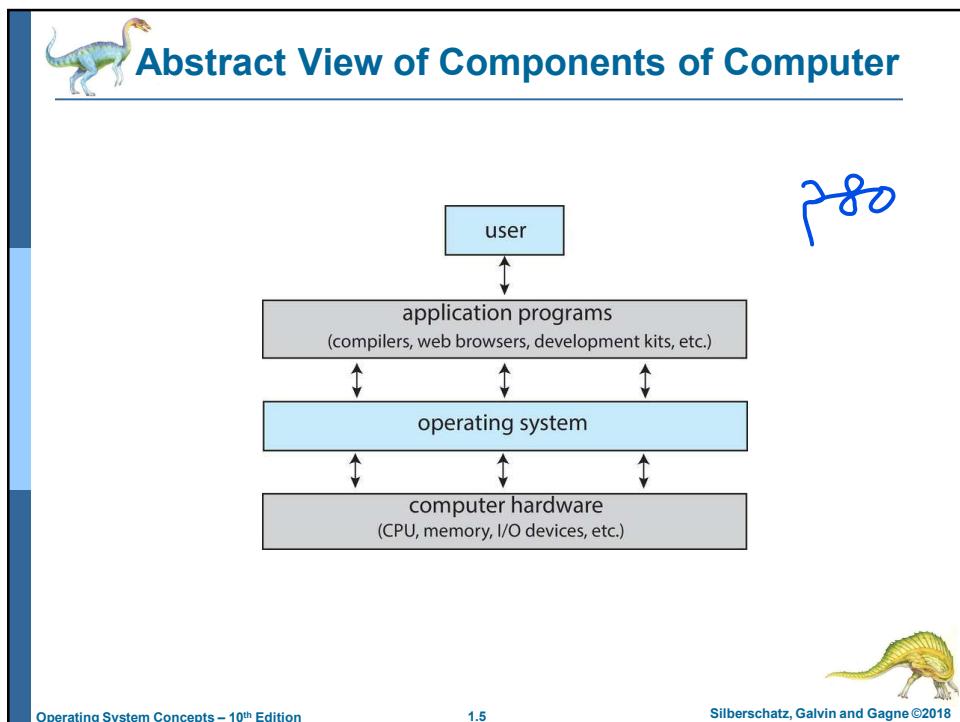
- A program that acts as an intermediary between a user of a computer and the computer hardware
- Operating system goals:
 1. Execute user programs and make solving user problems easier
 2. Make the computer system convenient to use
 3. Use the computer hardware in an efficient manner



Computer System Structure

- Computer system can be divided into four components:
 1. Hardware – provides basic computing resources
 - CPU, memory, I/O devices
 2. Operating system
 - Controls and coordinates use of hardware among various applications and users
 3. 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
 4. Users
 - People, machines, other computers



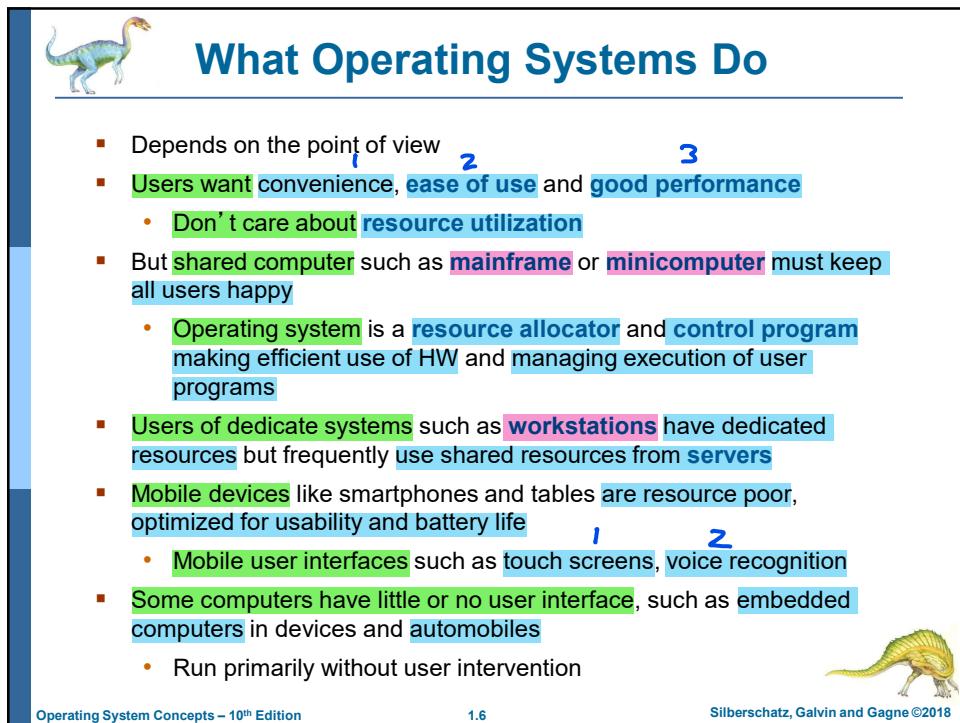


Operating System Concepts – 10th Edition

1.5

Silberschatz, Galvin and Gagne ©2018

5



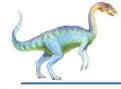
Operating System Concepts – 10th Edition

1.6

Silberschatz, Galvin and Gagne ©2018

6

3



Overview of Computer System Structure

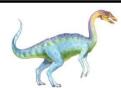


Operating System Concepts – 10th Edition

1.7

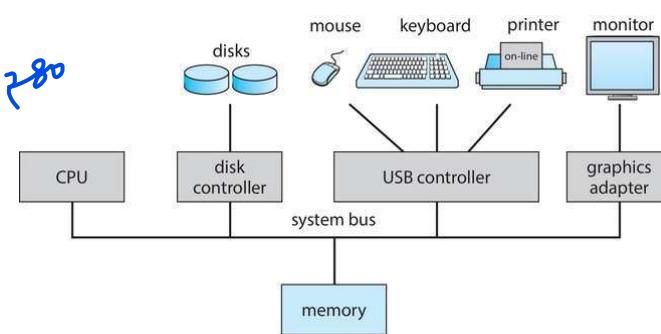
Silberschatz, Galvin and Gagne ©2018

7



Computer System Organization

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



Operating System Concepts – 10th Edition

1.8

Silberschatz, Galvin and Gagne ©2018

8

4



Computer-System Operation

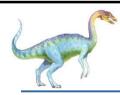
- I/O devices and the CPU can execute concurrently
 - Each device controller is in charge of a particular device type
 - Each device controller has a local buffer
 - Each device controller type has an operating system device driver to manage it
-
- CPU moves data from/to main memory to/from local buffers
 - 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



Operating-System Operations

- Bootstrap program – simple code to initialize the system, load the kernel
- Kernel loads
- 1 ▪ Starts system daemons (services provided outside of the kernel)
- 2 ▪ Kernel interrupt driven (hardware and software)
 - 1. ▪ Hardware interrupt by one of the devices
 - 2. ▪ Software interrupt (exception or trap):
 - Software error (e.g., division by zero)
 - Request for operating system service – system call
 - Other process problems include infinite loop, processes modifying each other or the operating system





Multiprogramming (Batch system)

- Single user cannot always keep CPU and I/O devices busy
- 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 job has to wait (for I/O for example), OS switches to another job

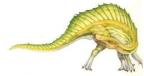
Jobs → code + data

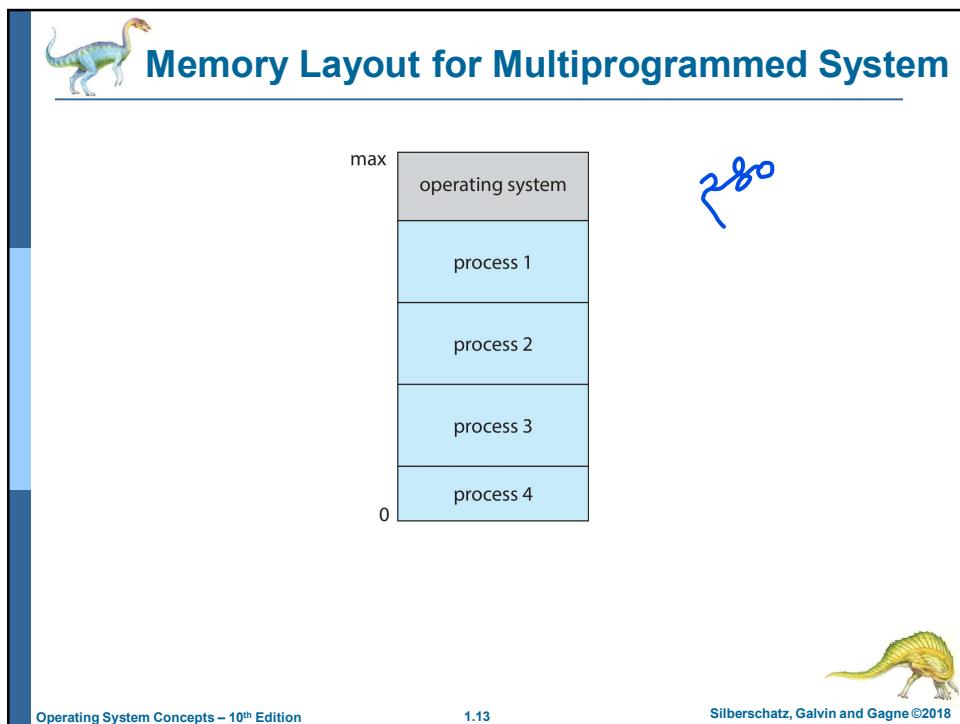


Multitasking (Timesharing)

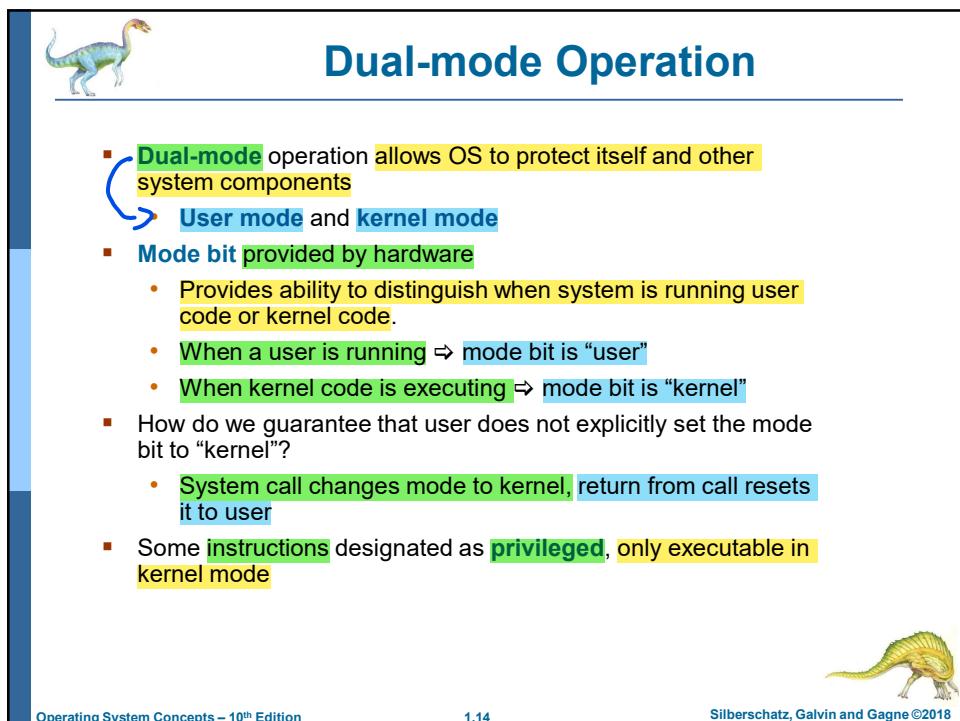


- A logical extension of Batch systems— the CPU switches jobs so frequently that users can interact with each job while it is running, creating interactive computing
 - Response time should be < 1 second
 - Each user has at least one program executing in memory ⇒ process
 - If several jobs ready to run 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

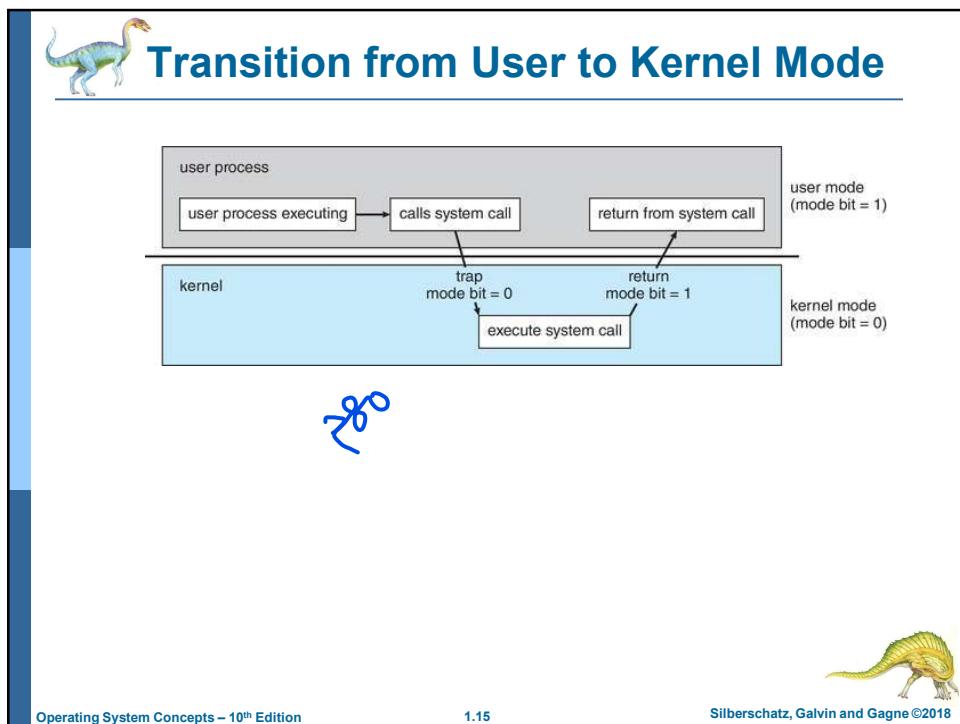




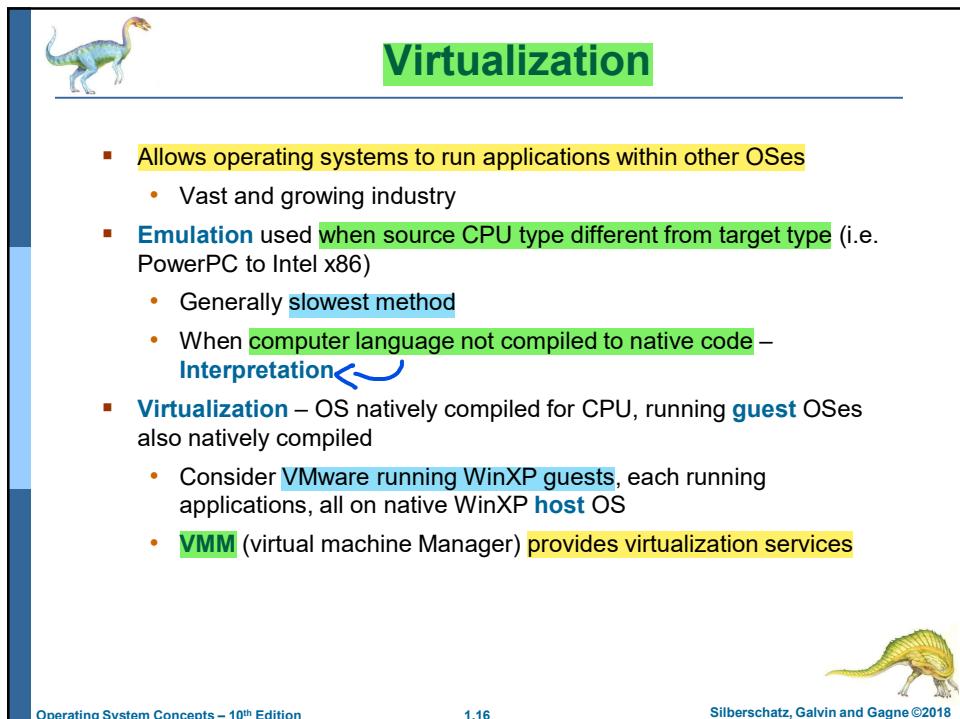
13



14



15



16

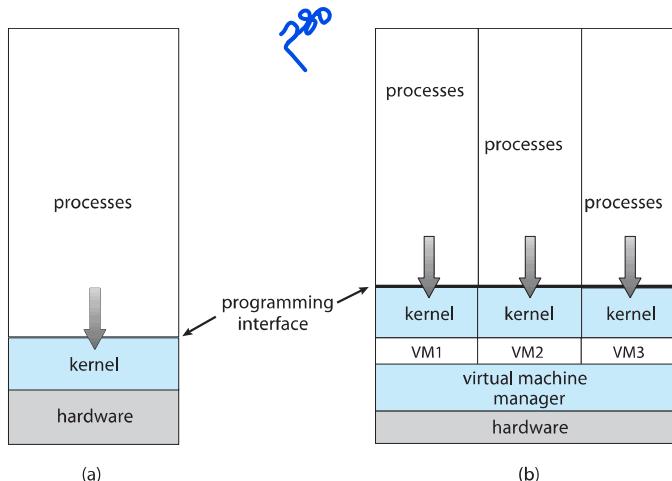


Virtualization (cont.)

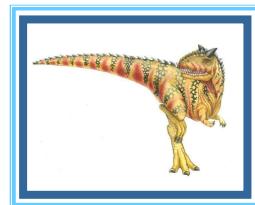
- **Use cases** involve laptops and desktops running multiple OSes for exploration or compatibility
 - Apple laptop running Mac OS X host, Windows as a guest
 - Developing apps for multiple OSes without having multiple systems
 - Quality assurance testing applications without having multiple systems
 - Executing and managing compute environments within data centers
- **VMM** can run natively, in which case they are also the host
 - There is no general-purpose host then (VMware ESX and Citrix XenServer)



Computing Environments - Virtualization



Chapter 1(Cont'd): Operating-System Services



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

19



Operating System Services

- Operating systems provide an environment for execution of programs and services to programs and users
- One set of operating-system services provides functions that are helpful to the user:
 1. **User interface** - Almost all operating systems have a user interface (**UI**).
 - Varies between **Command-Line (CLI)**, **Graphics User Interface (GUI)**, **touch-screen**, **Batch**
 2. **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)
 3. **I/O operations** - A running program may require I/O, which may involve a file or an I/O device
 4. **File-system manipulation** - The file system is of particular interest. Programs need to read and write files and directories, create and delete them, search them, list file information, permission management.



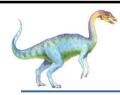
Operating System Concepts – 10th Edition

1.20

Silberschatz, Galvin and Gagne ©2018

20

10



Operating System Services (Cont.)

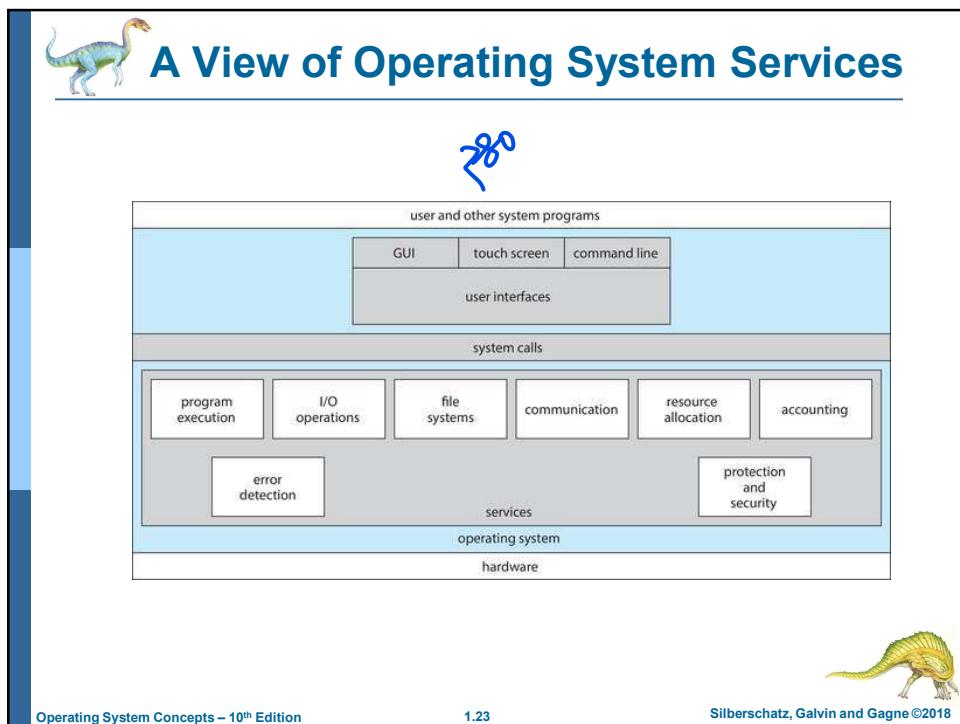
- One set of operating-system services provides functions that are helpful to the user (Cont.):
 - 5. **Communications** – Processes may exchange information, on the same computer or between computers over a network
 - ▶ Communications may be via shared memory or through message passing (packets moved by the OS)
 - 6. **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



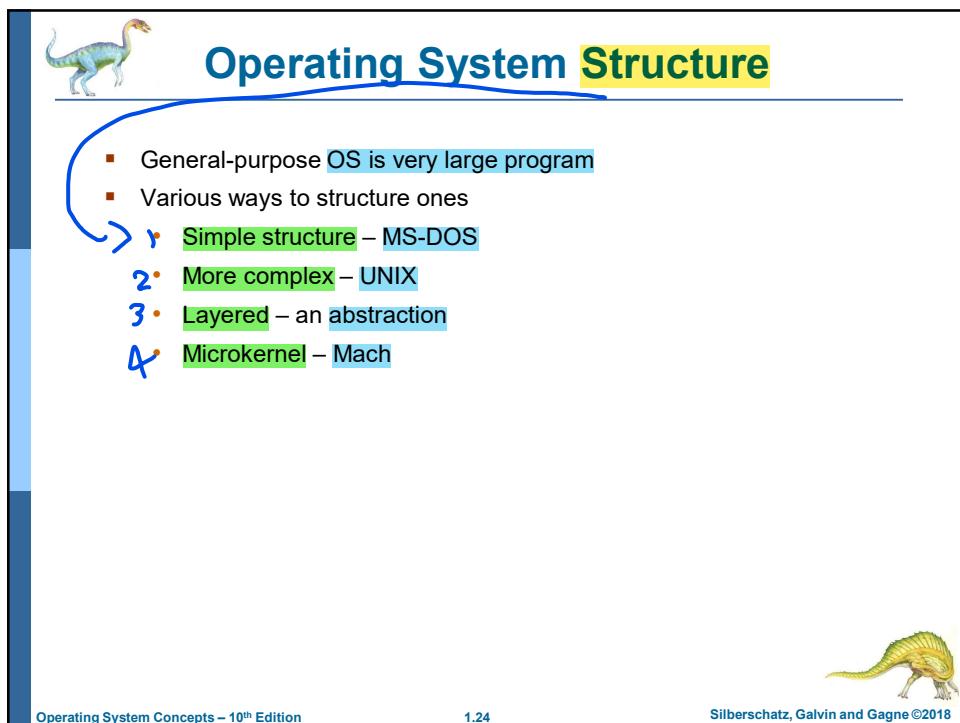
Operating System Services (Cont.)

- Another set of OS functions exists for ensuring the efficient operation of the system itself via resource sharing
 - 7. **Resource allocation** - When multiple users or multiple jobs running concurrently, resources must be allocated to each of them
 - ▶ Many types of resources - CPU cycles, main memory, file storage, I/O devices.
 - 8. **Logging** - To keep track of which users use how much and what kinds of computer resources
 - 9. **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





23



24



Monolithic Structure – Original UNIX

- **UNIX** – limited by hardware functionality, the original UNIX operating system had limited structuring.
- The **UNIX** OS consists of two separable **parts**
 - 1. Systems **programs**
 - 2. 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



Operating System Concepts – 10th Edition

1.25

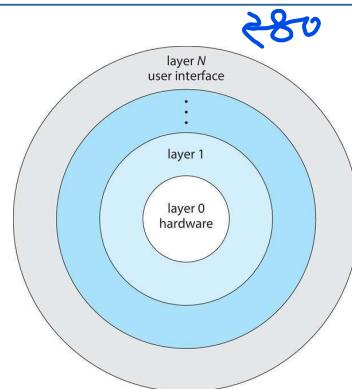
Silberschatz, Galvin and Gagne ©2018

25



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



Operating System Concepts – 10th Edition

1.26

Silberschatz, Galvin and Gagne ©2018

26

13



Microkernels

- Moves as much from the kernel into user space
- Mach is an example of microkernel
 - Mac OS X kernel (Darwin) partly based on Mach
- 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

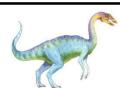


Operating System Concepts – 10th Edition

1.27

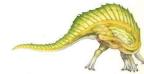
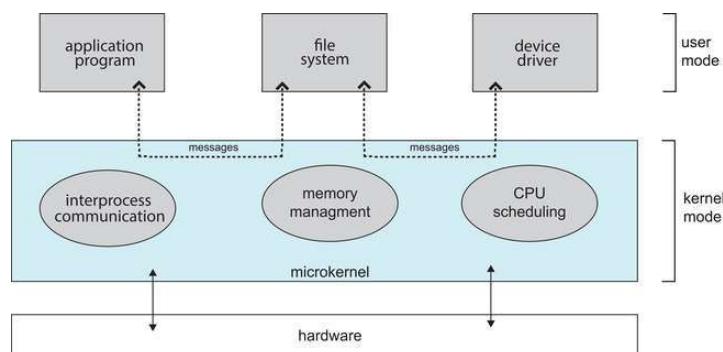
Silberschatz, Galvin and Gagne ©2018

27



Microkernel System Structure

280



Operating System Concepts – 10th Edition

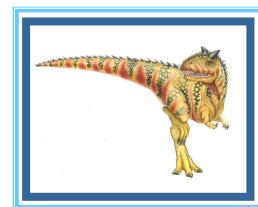
1.28

Silberschatz, Galvin and Gagne ©2018

28

14

Chapter 2: Processes & Threads



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

29



Outline

- Process Concept
- Process Scheduling
- Operations on Processes

Operating System Concepts – 10th Edition

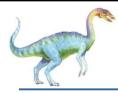
3.30

Silberschatz, Galvin and Gagne ©2018



30

15



Process Concept

- An operating system executes a variety of programs that run as a process.
- **Process** – a program in execution; process execution **must progress in sequential fashion**. No parallel execution of instructions of a single process
- Multiple parts
 - 1 The **program code**, also called **text section**
 - 2 Current activity including **program counter**, processor registers
 - 3 **Stack** containing **temporary data**
 - Function parameters, return addresses, local variables
 - 4 **Data section** containing **global variables**
 - 5 **Heap** containing **memory dynamically allocated** during run time



Process Concept (Cont.)

- Program is **passive** entity stored on disk (**executable file**); process is **active**
 - Program becomes process when an executable file is loaded into memory
- Execution of program started via **GUI** mouse clicks, **command line** entry of its name, etc.
- One **program can be several processes**
 - Consider multiple users executing the same program





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 processor
 - **Terminated:** The process has finished execution



Operating System Concepts – 10th Edition

3.33

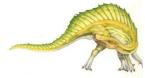
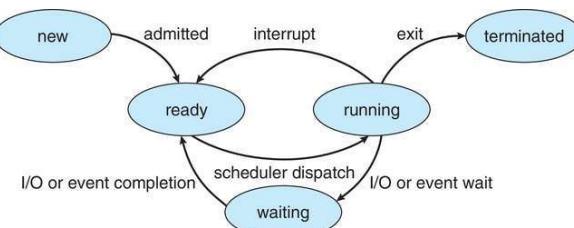
Silberschatz, Galvin and Gagne ©2018

33



Diagram of Process State

280

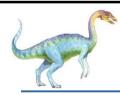


Operating System Concepts – 10th Edition

3.34

Silberschatz, Galvin and Gagne ©2018

34



Process Control Block (PCB)

Information associated with each process (also called **task control block**)

- 1 • Process state – running, waiting, etc.
- 2 • Program counter – location of instruction to next execute
- 3 • CPU registers – contents of all process-centric registers
- 4 • CPU scheduling information – priorities, scheduling queue pointers
- 5 • Memory-management information – memory allocated to the process
- 6 • Accounting information – CPU used, clock time elapsed since start, time limits
- 7 • I/O status information – I/O devices allocated to process, list of open files

process state
process number
program counter
registers
memory limits
list of open files
• • •



Operating System Concepts – 10th Edition

3.35

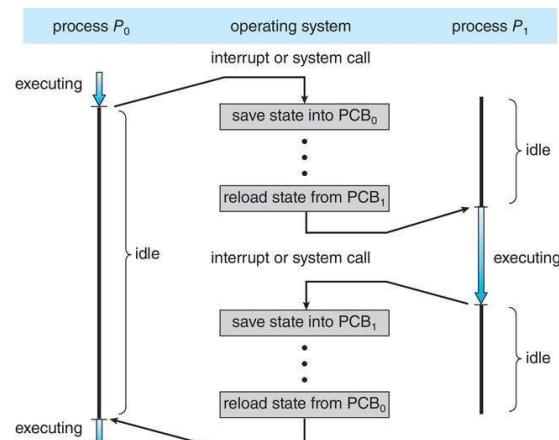
Silberschatz, Galvin and Gagne ©2018

35



CPU Switch From Process to Process

A **context switch** occurs when the CPU switches from one process to another.

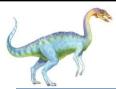


Operating System Concepts – 10th Edition

3.36

Silberschatz, Galvin and Gagne ©2018

36



Context Switch

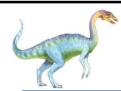
- When CPU switches to another process, the system must **save the state** of the old process and load the **saved state** for the new process via a **context switch**
- **Context** of a process **represented in the PCB**
- **Context-switch time** is **pure overhead**; the system does no useful work while switching
 - The **more complex** the OS and the PCB → the **longer** the context switch
- **Time dependent** on hardware support
 - Some hardware provides **multiple sets of registers per CPU** → **multiple contexts loaded at once**



Operations on Processes

- System must provide mechanisms for:
 - **Process creation**
 - **Process termination**





Process Creation

- Parent process create children processes, which, in turn create other processes, forming a tree of processes
- Generally, process identified and managed via a process identifier (pid)
- Resource sharing options
 1. Parent and children share all resources
 2. Children share subset of parent's resources
 3. Parent and child share no resources
- Execution options
 - Parent and children execute concurrently
 - Parent waits until children terminate



Operating System Concepts – 10th Edition

3.39

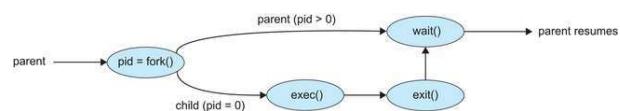
Silberschatz, Galvin and Gagne ©2018

39



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
 - Parent process calls `wait()` waiting for the child to terminate

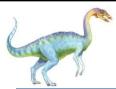


Operating System Concepts – 10th Edition

3.40

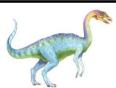
Silberschatz, Galvin and Gagne ©2018

40



Process Termination

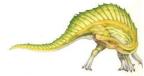
- Process executes last statement and then asks the operating system to delete it using the `exit()` system call.
 - Returns status data from child to parent (via `wait()`)
 - Process' resources are deallocated by operating system
- Parent may terminate the execution of children processes using the `abort()` system call. Some reasons for doing so:
 - Child has exceeded allocated resources
 - Task assigned to child is no longer required
 - The parent is exiting, and the operating systems does not allow a child to continue if its parent terminates



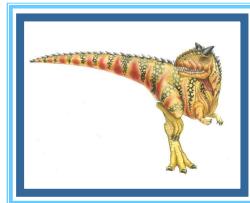
Process Termination

- Some operating systems do not allow child to exists if its parent has terminated. If a process terminates, then all its children must also be terminated.
 - **cascading termination.** All children, grandchildren, etc., are terminated.
 - The termination is initiated by the operating system.
- The parent process may wait for termination of a child process by using the `wait()` system call . The call returns status information and the pid of the terminated process

```
pid = wait(&status);
```
- If no parent waiting (did not invoke `wait()`) process is a **zombie**
- If parent terminated without invoking `wait()`, process is an **orphan**



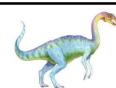
Chapter 2 (Cont'd): Threads



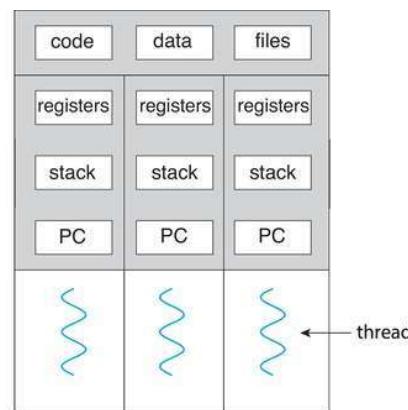
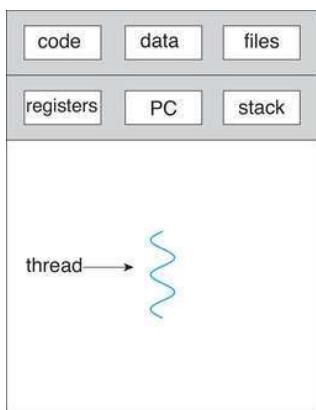
Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

43



Single and Multithreaded Processes



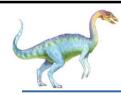
Operating System Concepts – 10th Edition

3.44

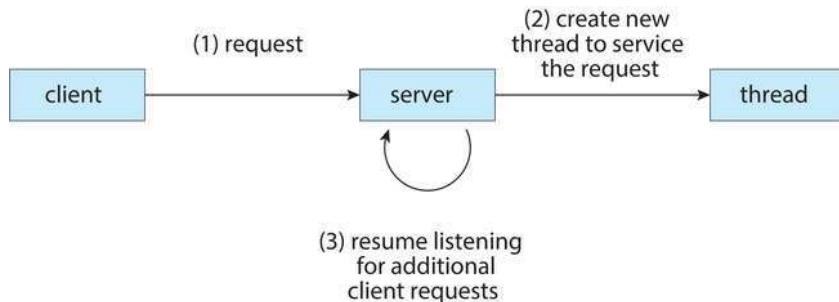
Silberschatz, Galvin and Gagne ©2018



44



Multithreaded Server Architecture



Operating System Concepts – 10th Edition

3.45

Silberschatz, Galvin and Gagne ©2018



45



Benefits

- **Responsiveness** – may allow continued execution if part of process is blocked, especially important for user interfaces
- **Resource Sharing** – threads share resources of process, easier than shared memory or message passing
- **Economy** – cheaper than process creation, thread switching lower overhead than context switching
- **Scalability** – process can take advantage of multicore architectures

Operating System Concepts – 10th Edition

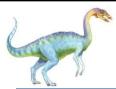
3.46

Silberschatz, Galvin and Gagne ©2018



46

23



Multicore Programming

- **Multicore** or **multiprocessor** systems puts pressure on programmers, challenges include:
 - Dividing activities
 - Balance
 - Data splitting
 - Data dependency
 - Testing and debugging
- **Parallelism** implies a system can perform more than one task simultaneously
- **Concurrency** supports more than one task making progress
 - Single processor / core, scheduler providing concurrency



Operating System Concepts – 10th Edition

3.47

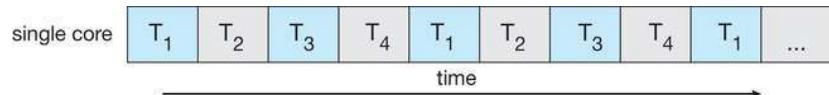
Silberschatz, Galvin and Gagne ©2018

47

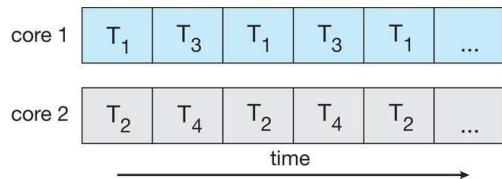


Concurrency vs. Parallelism

- Concurrent execution on single-core system:



- Parallelism on a multi-core system:

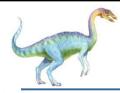


Operating System Concepts – 10th Edition

3.48

Silberschatz, Galvin and Gagne ©2018

48

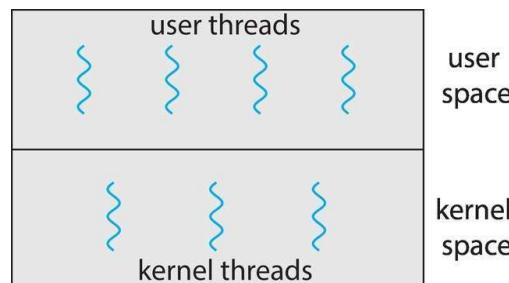


User Threads and Kernel Threads

- **User threads** - management done by user-level threads library
- Three primary thread libraries:
 - POSIX **Pthreads**
 - Windows threads
 - Java threads
- **Kernel threads** - Supported by the Kernel
- Examples – virtually all general-purpose operating systems, including:
 - Windows
 - Linux
 - Mac OS X
 - iOS
 - Android



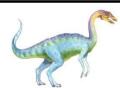
User and Kernel Threads





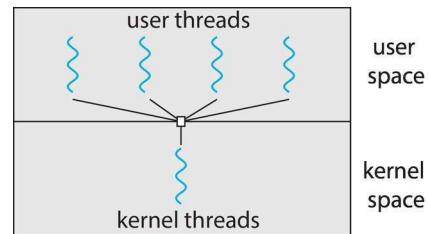
Multithreading Models

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



Many-to-One

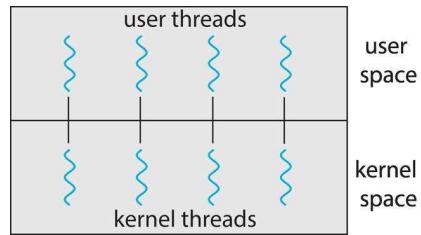
- Many user-level threads mapped to single kernel thread
- One thread blocking causes all to block
- Multiple threads may not run in parallel on multicore system because only one may be in kernel at a time
- Few systems currently use this model
- Examples:
 - Solaris Green Threads
 - GNU Portable Threads





One-to-One

- Each user-level thread maps to kernel thread
- Creating a user-level thread creates a kernel thread
- More concurrency than many-to-one
- Number of threads per process sometimes restricted due to overhead
- Examples
 - Windows
 - Linux



Operating System Concepts – 10th Edition

3.53

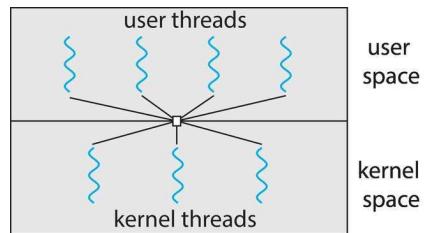
Silberschatz, Galvin and Gagne ©2018

53



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
- Windows with the *ThreadFiber* package
- Otherwise not very common

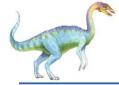


Operating System Concepts – 10th Edition

3.54

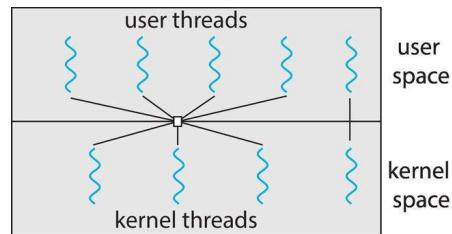
Silberschatz, Galvin and Gagne ©2018

54



Two-level Model

- Similar to M:M, except that it allows a user thread to be **bound** to kernel thread



Thread Libraries

- **Thread library** provides programmer with API for creating and managing threads
- Two primary ways of implementing
 - Library entirely in user space
 - Kernel-level library supported by the OS





Thread Cancellation

- Terminating a thread before it has finished
- Thread to be canceled is **target thread**
- Two general approaches:
 - **Asynchronous cancellation** terminates the target thread immediately
 - **Deferred cancellation** allows the target thread to periodically check if it should be cancelled
- Pthread code to create and cancel a thread:

```
pthread_t tid;  
  
/* create the thread */  
pthread_create(&tid, 0, worker, NULL);  
  
.  
.  
.  
  
/* cancel the thread */  
pthread_cancel(tid);  
  
/* wait for the thread to terminate */  
pthread_join(tid,NULL);
```

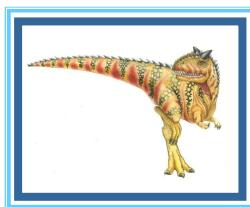


Thread Cancellation (Cont.)

- Invoking thread cancellation requests cancellation, but actual cancellation depends on thread state
 - Mode State Type
 - Off Disabled –
 - Deferred Enabled Deferred
 - Asynchronous Enabled Asynchronous
- If thread has cancellation disabled, cancellation remains pending until thread enables it
 - Default type is deferred
 - Cancellation only occurs when thread reaches **cancellation point**
 - i.e., `pthread_testcancel()`
 - Then **cleanup handler** is invoked
 - On Linux systems, thread cancellation is handled through signals



Chapter 3: CPU Scheduling



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

59



Outline

- Basic Concepts
- Scheduling Criteria
- Scheduling Algorithms
- Multi-Processor Scheduling
- Real-Time CPU Scheduling

Operating System Concepts – 10th Edition

5.60

Silberschatz, Galvin and Gagne ©2018



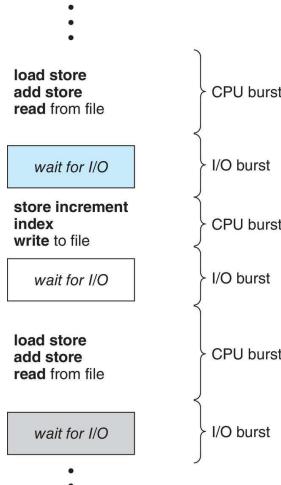
60

30



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** followed by **I/O burst**
- CPU burst distribution is of main concern



Operating System Concepts – 10th Edition

5.61

Silberschatz, Galvin and Gagne ©2018

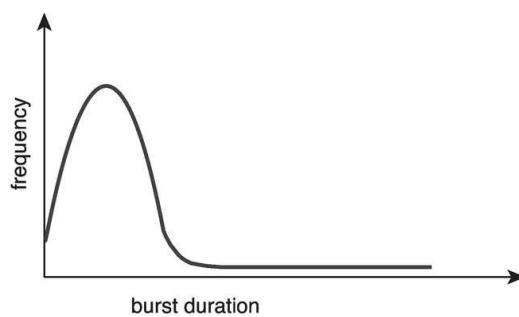
61



Histogram of CPU-burst Times

Large number of short bursts

Small number of longer bursts



Operating System Concepts – 10th Edition

5.62

Silberschatz, Galvin and Gagne ©2018

62

31



CPU Scheduler

- The **CPU scheduler** selects from among the processes in ready queue, and allocates a CPU core to one of them
 - Queue may be ordered in various ways
- 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
- For situations 1 and 4, there is no choice in terms of scheduling. A new process (if one exists in the ready queue) must be selected for execution.
- For situations 2 and 3, however, there is a choice.



Preemptive and Nonpreemptive Scheduling

- When scheduling takes place only under circumstances 1 and 4, the scheduling scheme is **nonpreemptive**.
- Otherwise, it is **preemptive**.
- Under Nonpreemptive scheduling, once the CPU has been allocated to a process, the process keeps the CPU until it releases it either by terminating or by switching to the waiting state.
- Virtually all modern operating systems including Windows, MacOS, Linux, and UNIX use preemptive scheduling algorithms.





Preemptive Scheduling and Race Conditions

- Preemptive scheduling can result in race conditions when data are shared among several processes.
- Consider the case of two processes that share data. While one process is updating the data, it is preempted so that the second process can run. The second process then tries to read the data, which are in an inconsistent state.
- This issue will be explored in detail in Chapter 6.



Operating System Concepts – 10th Edition

5.65

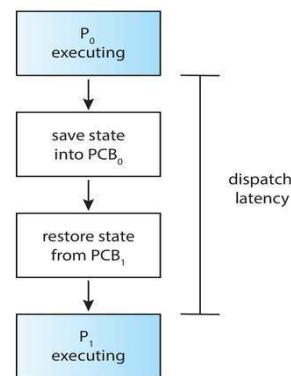
Silberschatz, Galvin and Gagne ©2018

65



Dispatcher

- Dispatcher module gives control of the CPU to the process selected by the CPU 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



Operating System Concepts – 10th Edition

5.66

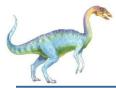
Silberschatz, Galvin and Gagne ©2018

66



Scheduling Criteria

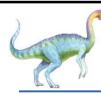
- **CPU utilization** – keep the CPU as busy as possible
- **Throughput** – # of processes that complete their execution per time unit
- **Turnaround time** – amount of time to execute a particular process
- **Waiting time** – amount of time a process has been waiting in the ready queue
- **Response time** – amount of time it takes from when a request was submitted until the first response is produced.



Scheduling Algorithm Optimization Criteria

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





First-Come, First-Served (FCFS) Scheduling

First Come First Serve (FCFS) Scheduling

- Jobs are executed on first come, first serve basis.
- Easy to understand and implement.
- Poor in performance as average wait time is high.

PROCESS	BURST TIME
P1	21
P2	3
P3	6
P4	2

The GANTT chart for FCFS will be,



PROCESS	WAITING TIME
P1	0
P2	21
P3	24
P4	30

The average waiting time will be = $(0+21+24+30) / 4 = 18.75 \text{ ms}$



Page No: 3



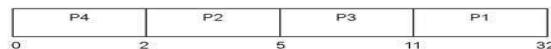
Shortest-Job-First (SJF) Scheduling

Shortest-Job-First (SJF) Scheduling – Non Preemptive

- Best approach to minimize waiting time.
- Actual time taken by the process is already known to processor.
- Shortest process is executed first.

PROCESS	BURST TIME
P1	21
P2	3
P3	6
P4	2

The GANTT chart for SJF-Non Preemptive will be,



PROCESS	WAITING TIME
P1	11
P2	2
P3	5
P4	0

The average waiting time will be = $(11+2+5+0) / 4 = 4.5 \text{ ms}$



Page No: 4

 Example of SJF(Preemptive)

Shortest-Job-First (SJF) Scheduling – Preemptive

In Preemptive Shortest Job First Scheduling, jobs are put into ready queue as they arrive, but as a process with short burst time arrives, the existing process is preempted.

PROCESS	BURST TIME	ARRIVAL TIME
P1	21	0
P2	3	1
P3	6	2
P4	2	3

The GANTT chart for SJF-PREEMPTIVE will be,

P1	P2	P4	P3	P1
0	1	4	6	12
32				

PROCESS	BURST TIME (BT)	ARRIVAL TIME (AT)	COMPLETION TIME (CT)	TURNAROUND TIME(TAT) TAT = CT-AT	WAITING TIME (WT) WT=TAT-BT
P1	21	0	32	32	11
P2	3	1	4	3	0
P3	6	2	12	10	4
P4	2	3	6	3	1
TOTAL WAITING TIME				16	

The average waiting time will be $(11+0+4+1) / 4 = 4 \text{ ms}$



Page No: 5
Silberschatz, Galvin and Gagne ©2018

Operating System Concepts – 10th Edition 5.71

71

 Round Robin (RR)

- A fixed time is allocated to each process, called **quantum** for each execution
- Only a process is executed for given time period that process is preempted and other process executes for given time period
- Context switching is used to save states of preempted processes.
- Assume time quantum as 5ms

Process	Burst time
P1	21
P2	3
P3	6
P4	2

Gant Chart

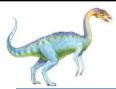
P1	P2	P3	P4	P1	P3	P1
0	5	8	13	15	20	21
32						

Process	Waiting Time = Completion Time – Burst time
P1	$= (32-21) \rightarrow 11$
P2	$= (8-3) \rightarrow 5$
P3	$= (21-6) \rightarrow 15$
P4	$= (15-2) \rightarrow 13$
Total Waiting Time	=44
Average Waiting Time	$=44/4 \rightarrow 11 \text{ ms}$



Operating System Concepts – 10th Edition 5.72
Silberschatz, Galvin and Gagne ©2018

72



Priority Scheduling

- Priority is assigned for each process.
- Process with highest priority is executed first and so on.
- Process with same priority are executed in FCFS manner.
- Priority can be decided on memory requirements, time requirements or any other resource requirements

Proces	Burst Time	Priority
P1	21	2
P2	3	1 (Highest)
P3	6	4 (Lowest)
P4	2	3



Waiting Time

Total Wt = 53

AWT = $53/4 = 18.75\text{ms}$

Process	Waiting Time
P1	3
P2	0
P3	26
P4	24

Operating System Concepts – 10th Edition

5.73

Silberschatz, Galvin and Gagne ©2018



73



Multilevel Queue

- The ready queue consists of multiple queues
- Multilevel queue scheduler defined by the following parameters:
 - Number of queues
 - Scheduling algorithms for each queue
 - Method used to determine which queue a process will enter when that process needs service
 - Scheduling among the queues

Operating System Concepts – 10th Edition

5.74

Silberschatz, Galvin and Gagne ©2018

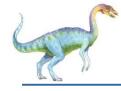
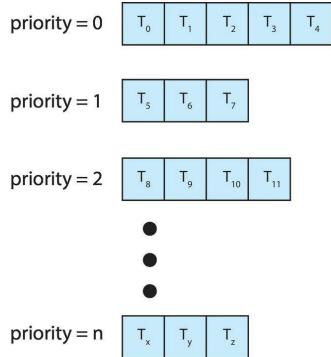


74



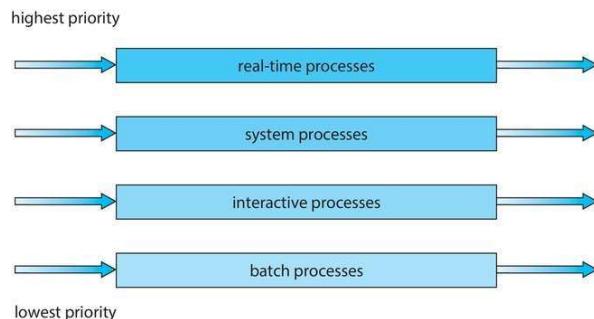
Multilevel Queue

- With priority scheduling, have separate queues for each priority.
- Schedule the process in the highest-priority queue!



Multilevel Queue

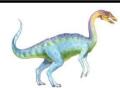
- Prioritization based upon process type





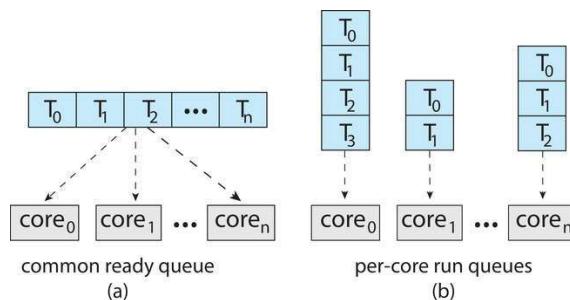
Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available
- Multiprocess may be any one of the following architectures:
 - Multicore CPUs
 - Multithreaded cores
 - NUMA systems
 - Heterogeneous multiprocessing



Multiple-Processor Scheduling

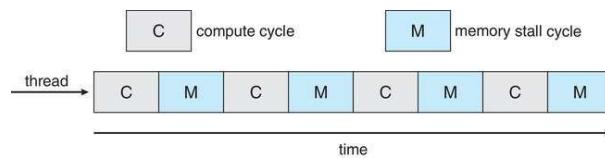
- Symmetric multiprocessing (SMP) is where each processor is self scheduling.
- All threads may be in a common ready queue (a)
- Each processor may have its own private queue of threads (b)





Multicore Processors

- Recent trend to place multiple processor cores on same physical chip
- Faster and consumes less power
- Multiple threads per core also growing
 - Takes advantage of memory stall to make progress on another thread while memory retrieve happens
- Figure



Operating System Concepts – 10th Edition

5.79

Silberschatz, Galvin and Gagne ©2018



79



Real-Time CPU Scheduling

- Can present obvious challenges
- **Soft real-time systems** – Critical real-time tasks have the highest priority, but no guarantee as to when tasks will be scheduled
- **Hard real-time systems – task must be serviced by its deadline**

Operating System Concepts – 10th Edition

5.80

Silberschatz, Galvin and Gagne ©2018



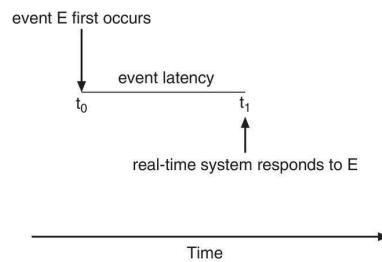
80

40



Real-Time CPU Scheduling

- Event latency – the amount of time that elapses from when an event occurs to when it is serviced.
- Two types of latencies affect performance
 1. **Interrupt latency** – time from arrival of interrupt to start of routine that services interrupt
 2. **Dispatch latency** – time for scheduler to take current process off CPU and switch to another



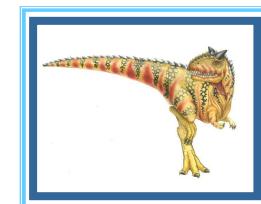
Operating System Concepts – 10th Edition

5.81

Silberschatz, Galvin and Gagne ©2018

81

Chapter 4: Synchronization & Deadlocks



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

82

41



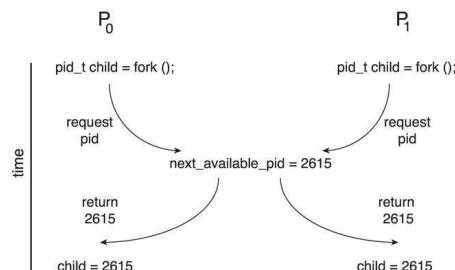
Outline

- Background
- The Critical-Section Problem
- Mutex Locks
- Semaphores
- Monitors



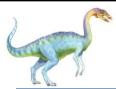
Race Condition

- Processes P_0 and P_1 are creating child processes using the `fork()` system call
- Race condition on kernel variable `next_available_pid` which represents the next available process identifier (pid)



- Unless there is a mechanism to prevent P_0 and P_1 from accessing the variable `next_available_pid` the same pid could be assigned to two different processes!





Critical Section Problem

- Consider system of n processes $\{p_0, p_1, \dots, p_{n-1}\}$
- Each process has **critical section** segment of code
 - Process may be changing common variables, updating table, writing file, etc.
 - When one process in critical section, no other may be in its critical section
- **Critical section problem** is to design protocol to solve this
- Each process must ask permission to enter critical section in **entry section**, may follow critical section with **exit section**, then **remainder section**



Critical-Section Problem (Cont.)

Requirements for solution to critical-section problem

1. **Mutual Exclusion** - If process P_i is executing in its critical section, then no other processes can be executing in their critical sections
2. **Progress** - If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the process 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





Mutex Locks

- Previous solutions are complicated and generally inaccessible to application programmers
- OS designers build software tools to solve critical section problem
- Simplest is mutex lock
 - Boolean variable indicating if lock is available or not
- Protect a critical section by
 - First **acquire()** a lock
 - Then **release()** the lock
- Calls to **acquire()** and **release()** must be **atomic**
 - Usually implemented via hardware atomic instructions such as compare-and-swap.
- But this solution requires **busy waiting**
 - This lock therefore called a **spinlock**



Semaphore

- Synchronization tool that provides more sophisticated ways (than Mutex locks) for processes to synchronize their activities.
- Semaphore **S** – integer variable
- Can only be accessed via two indivisible (atomic) operations
 - **wait()** and **signal()**
 - Originally called **P()** and **V()**
- Definition of the **wait()** operation

```
wait(S) {  
    while (S <= 0)  
        ; // busy wait  
    S--;  
}
```
- Definition of the **signal()** operation

```
signal(S) {  
    S++;  
}
```



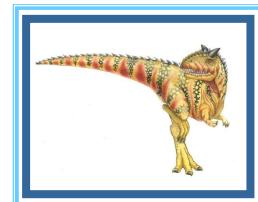


Semaphore (Cont.)

- **Counting semaphore** – integer value can range over an unrestricted domain
- **Binary semaphore** – integer value can range only between 0 and 1
 - Same as a **mutex lock**
- Can implement a counting semaphore **S** as a binary semaphore
- With semaphores we can solve various synchronization problems



Chapter 4 (Cont'd): Synchronization Examples





Classical Problems of Synchronization

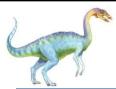
- Classical problems used to test newly-proposed synchronization schemes
 - Bounded-Buffer Problem
 - Readers and Writers Problem
 - Dining-Philosophers Problem



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





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
- Several variations of how readers and writers are considered – all involve some form of priorities

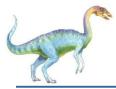


Operating System Concepts – 10th Edition

6.93

Silberschatz, Galvin and Gagne ©2018

93



Readers-Writers Problem (Cont.)

- Shared Data
 - Data set
 - Semaphore `rw_mutex` initialized to 1
 - Semaphore `mutex` initialized to 1
 - Integer `read_count` initialized to 0



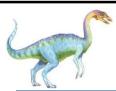
Operating System Concepts – 10th Edition

6.94

Silberschatz, Galvin and Gagne ©2018

94

47



Dining-Philosophers Problem

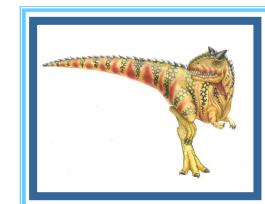
- N philosophers sit at a round table with a bowel of rice in the middle.



- They spend their lives alternating thinking and eating.
- They do not interact with their neighbors.
- Occasionally try to pick up 2 chopsticks (one at a time) to eat from bowl
 - Need both to eat, then release both when done
- In the case of 5 philosophers, the shared data
 - ▶ Bowl of rice (data set)
 - ▶ Semaphore chopstick [5] initialized to 1



Chapter 4 (Cont'd): Deadlocks





System Model

- System consists of resources
- Resource types R_1, R_2, \dots, 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

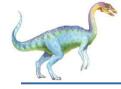


Deadlock Characterization

Deadlock can arise if four conditions hold simultaneously.

- **Mutual exclusion:** only one thread at a time can use a resource
- **Hold and wait:** a thread holding at least one resource is waiting to acquire additional resources held by other threads
- **No preemption:** a resource can be released only voluntarily by the thread holding it, after that thread has completed its task
- **Circular wait:** there exists a set $\{T_0, T_1, \dots, T_n\}$ of waiting threads such that T_0 is waiting for a resource that is held by T_1 , T_1 is waiting for a resource that is held by T_2, \dots, T_{n-1} is waiting for a resource that is held by T_n , and T_n is waiting for a resource that is held by T_0 .





Resource-Allocation Graph

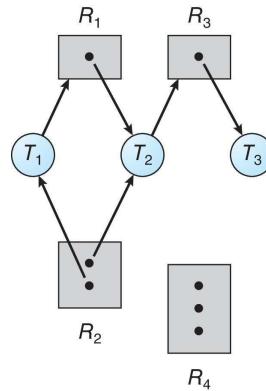
A set of vertices V and a set of edges E .

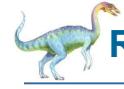
- V is partitioned into two types:
 - $T = \{T_1, T_2, \dots, T_n\}$, the set consisting of all the threads in the system.
 - $R = \{R_1, R_2, \dots, R_m\}$, the set consisting of all resource types in the system
- **request edge** – directed edge $T_i \rightarrow R_j$
- **assignment edge** – directed edge $R_j \rightarrow T_i$



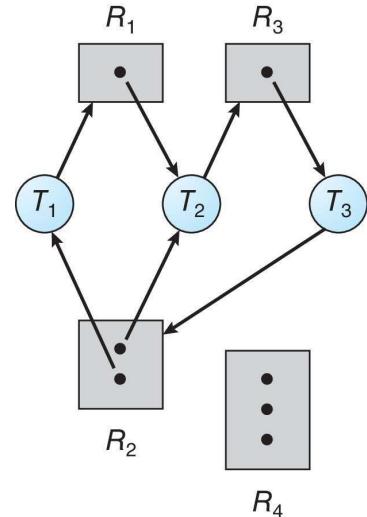
Resource Allocation Graph Example

- One instance of R_1
- Two instances of R_2
- One instance of R_3
- Three instance of R_4
- T_1 holds one instance of R_2 and is waiting for an instance of R_1
- T_2 holds one instance of R_1 , one instance of R_2 , and is waiting for an instance of R_3
- T_3 is holds one instance of R_3





Resource Allocation Graph with a Deadlock

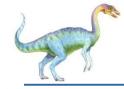


Operating System Concepts – 10th Edition

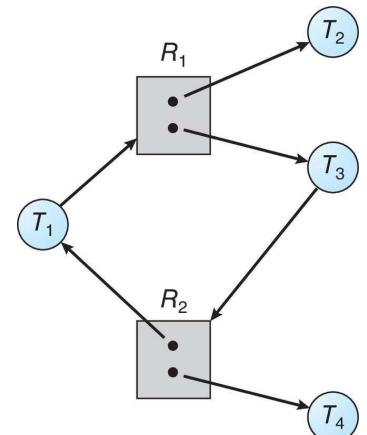
6.101

Silberschatz, Galvin and Gagne ©2018

101



Graph with a Cycle But no Deadlock

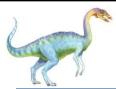


Operating System Concepts – 10th Edition

6.102

Silberschatz, Galvin and Gagne ©2018

102



Basic Facts

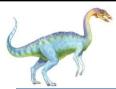
- If graph contains no cycles \Rightarrow no deadlock
- If graph contains a cycle \Rightarrow
 - if only one instance per resource type, then deadlock
 - if several instances per resource type, possibility of deadlock



Methods for Handling Deadlocks

- Ensure that the system will **never** enter a deadlock state:
 - Deadlock prevention
 - Deadlock avoidance
- Allow the system to enter a deadlock state and then recover
- Ignore the problem and pretend that deadlocks never occur in the system.





Deadlock Prevention

Invalidate one of the four necessary conditions for deadlock:

- **Mutual Exclusion** – not required for sharable resources (e.g., read-only files); must hold for non-sharable resources
- **Hold and Wait** – must guarantee that whenever a thread requests a resource, it does not hold any other resources
 - Require threads to request and be allocated all its resources before it begins execution or allow thread to request resources only when the thread has none allocated to it.
 - 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 thread is waiting
- Thread 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 thread requests resources in an increasing order of enumeration





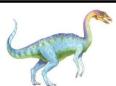
Circular Wait

- Invalidating the circular wait condition is most common.
- Simply assign each resource (i.e., mutex locks) a unique number.
- Resources must be acquired in order.
- If:

```
first_mutex = 1  
second_mutex = 5
```

code for **thread_two** could not be written as follows:

```
/* thread.one runs in this function */  
void *do_work_one(void *param)  
{  
    pthread_mutex_lock(&first_mutex);  
    pthread_mutex_lock(&second_mutex);  
    /*  
     * Do some work  
     */  
    pthread_mutex_unlock(&second_mutex);  
    pthread_mutex_unlock(&first_mutex);  
  
    pthread_exit(0);  
}  
  
/* thread.two runs in this function */  
void *do_work_two(void *param)  
{  
    pthread_mutex_lock(&second_mutex);  
    pthread_mutex_lock(&first_mutex);  
    /*  
     * Do some work  
     */  
    pthread_mutex_unlock(&first_mutex);  
    pthread_mutex_unlock(&second_mutex);  
  
    pthread_exit(0);  
}
```

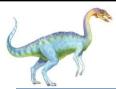


Deadlock Avoidance

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

- Simplest and most useful model requires that each thread 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 thread 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 sequence $\langle T_1, T_2, \dots, T_n \rangle$ of ALL the threads in the systems such that for each T_i , the resources that T_i can still request can be satisfied by currently available resources + resources held by all the T_j , with $j < i$
- That is:
 - If T_i resource needs are not immediately available, then T_i can wait until all T_j have finished
 - When T_j is finished, T_i can obtain needed resources, execute, return allocated resources, and terminate
 - When T_i terminates, T_{i+1} can obtain its needed resources, and so on



Operating System Concepts – 10th Edition

6.109

Silberschatz, Galvin and Gagne ©2018

109



Basic Facts

- If a system is in safe state \Rightarrow no deadlocks
- If a system is in unsafe state \Rightarrow possibility of deadlock
- Avoidance \Rightarrow ensure that a system will never enter an unsafe state.

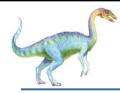


Operating System Concepts – 10th Edition

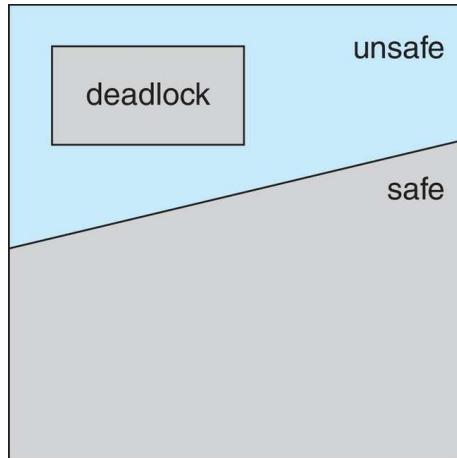
6.110

Silberschatz, Galvin and Gagne ©2018

110



Safe, Unsafe, Deadlock State



Operating System Concepts – 10th Edition

6.111

Silberschatz, Galvin and Gagne ©2018

111



Avoidance Algorithms

- Single instance of a resource type
 - Use a resource-allocation graph

- Multiple instances of a resource type
 - Use the Banker's Algorithm



Operating System Concepts – 10th Edition

6.112

Silberschatz, Galvin and Gagne ©2018

112



Banker's Algorithm

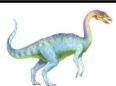
- Multiple instances of resources
- Each thread must a priori claim maximum use
- When a thread requests a resource, it may have to wait
- When a thread gets all its resources it must return them in a finite amount of time

Operating System Concepts – 10th Edition

6.113

Silberschatz, Galvin and Gagne ©2018

113



Example :

Banker's Algorithm

• Considering a system with five processes P0 through P4 and three resources types A, B, C. Resource type A has 5 instances, B has 10 instances and C has 7 instances. Suppose at time t_0 following snapshot of the system has been taken:

Process	Allocation			Max			Available		
	A	B	C	A	B	C	A	B	C
P0	1	0	0	5	7	3	3	3	2
P1	0	2	0	2	3	2			
P2	0	3	2	0	9	2			
P3	1	2	1	2	2	2			
P4	0	0	2	3	4	3			

By applying banker's algorithm

• What will be the content of need matrix?

• Is the system in safe state? If yes, then what is the safe sequence?

Step 1:

i) Need = Max – Alloc

Step 2:

• Work = Available = [3, 3, 2]

• Finish (i)= [F, F, F, F, F]

Step 3:

→ P0	4	7	3
P1	2	1	2
P2	0	6	0
P3	1	0	1
P4	3	4	2

Find the process where need<=work	If so, work=work+allocation	Finish (i) = True
P1	→[3,3,2] + [0,2,0] →[3, 5, 2]	[F, T, F, F, F]
P3	→[3,5,2] + [1, 2, 1] → [4,7,3]	[F, T, F, T, F]
P4	→[4,7,3] + [0,0,2] → [4,7,5]	[F, T, F, T, T]
P0	→[4,7,5] + [1,0,0] → [5,7,5]	[T, T, F, T, T]
P2	→[5,7,5] + [0,3,2] → [5,10,7]	[T, T, T, T, T]



ii) System is in safe state, and the safe sequence is:

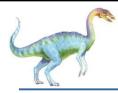
P1 → P3 → P4 → P0 → P2

Operating System Concepts – 10th Edition

6.114

Silberschatz, Galvin and Gagne ©2018

114



Deadlock Detection

- Allow system to enter deadlock state
- Detection algorithm
- Recovery scheme



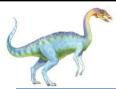
Example of Detection Algorithm

- Five threads T_0 through T_4 ; three resource types A (7 instances), B (2 instances), and C (6 instances)
- Snapshot at time T_0 :

	<u>Allocation</u>	<u>Request</u>	<u>Available</u>
	A B C	A B C	A B C
T_0	0 1 0	0 0 0	0 0 0
T_1	2 0 0	2 0 2	
T_2	3 0 3	0 0 0	
T_3	2 1 1	1 0 0	
T_4	0 0 2	0 0 2	

- Sequence $\langle T_0, T_2, T_3, T_1, T_4 \rangle$ will result in $Finish[i] = \text{true}$ for all i





Example (Cont.)

- T_2 requests an additional instance of type C

Request

	A	B	C
T_0	0	0	0
T_1	2	0	2
T_2	0	0	1
T_3	1	0	0
T_4	0	0	2

- State of system?

- Can reclaim resources held by thread T_0 , but insufficient resources to fulfill other processes' requests
- Deadlock exists, consisting of processes T_1 , T_2 , T_3 , and T_4



Recovery from Deadlock: Process Termination

- Abort all deadlocked threads
- Abort one process at a time until the deadlock cycle is eliminated
- In which order should we choose to abort?
 1. Priority of the thread
 2. How long has the thread computed, and how much longer to completion
 3. Resources that the thread has used
 4. Resources that the thread needs to complete
 5. How many threads will need to be terminated
 6. Is the thread interactive or batch?





Recovery from Deadlock: Resource Preemption

- **Selecting a victim** – minimize cost
- **Rollback** – return to some safe state, restart the thread for that state
- **Starvation** – same thread may always be picked as victim, include number of rollback in cost factor



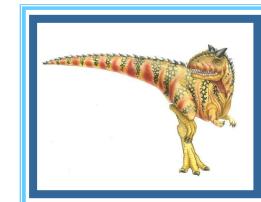
Operating System Concepts – 10th Edition

6.119

Silberschatz, Galvin and Gagne ©2018

119

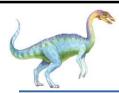
Chapter 5: Memory Management (Main Memory)



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

120



Chapter 5: Memory Management

- Background
- Contiguous Memory Allocation
- Paging
- Structure of the Page Table
- Swapping



Operating System Concepts – 10th Edition

9.121

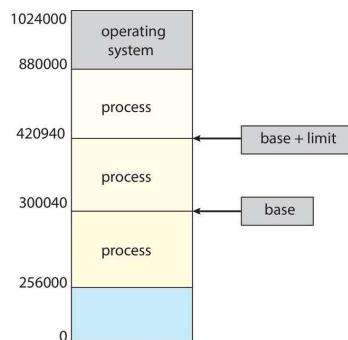
Silberschatz, Galvin and Gagne ©2018

121



Protection

- Need to ensure that a process can access only those addresses in its address space.
- We can provide this protection by using a pair of **base** and **limit registers** define the logical address space of a process

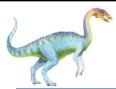


Operating System Concepts – 10th Edition

9.122

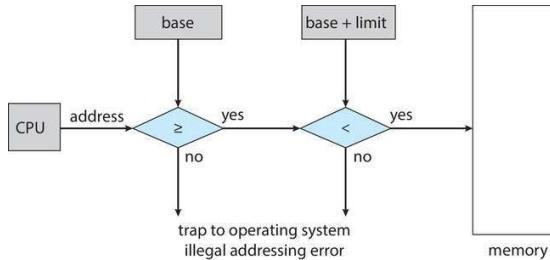
Silberschatz, Galvin and Gagne ©2018

122



Hardware Address Protection

- CPU must check every memory access generated in user mode to be sure it is between base and limit for that user



- the instructions to loading the base and limit registers are privileged



Address Binding

- Programs on disk, ready to be brought into memory to execute from an **input queue**
 - Without support, must be loaded into address 0000
- Inconvenient to have first user process physical address always at 0000
 - How can it not be?
- Addresses represented in different ways at different stages of a program's life
 - Source code addresses usually symbolic
 - Compiled code addresses **bind** to relocatable addresses
 - i.e., "14 bytes from beginning of this module"
 - Linker or loader will bind relocatable addresses to absolute addresses
 - i.e., 74014
 - Each binding maps one address space to another





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)

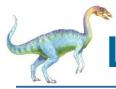


Operating System Concepts – 10th Edition

9.125

Silberschatz, Galvin and Gagne ©2018

125



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
- **Logical address space** is the set of all logical addresses generated by a program
- **Physical address space** is the set of all physical addresses generated by a program

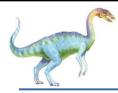


Operating System Concepts – 10th Edition

9.126

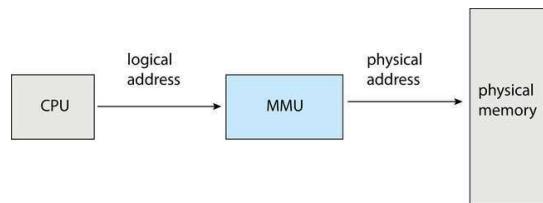
Silberschatz, Galvin and Gagne ©2018

126



Memory-Management Unit (MMU)

- Hardware device that at run time maps virtual to physical address



- Many methods possible, covered in the rest of this chapter

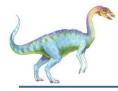


Operating System Concepts – 10th Edition

9.127

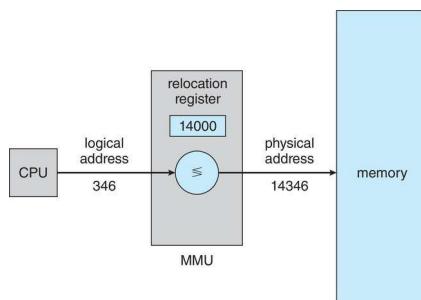
Silberschatz, Galvin and Gagne ©2018

127



Memory-Management Unit (Cont.)

- Consider simple scheme, which is a generalization of the base-register scheme.
- The base register now called **relocation register**
- The value in the relocation register is added to every address generated by a user process at the time it is sent to memory

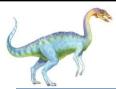


Operating System Concepts – 10th Edition

9.128

Silberschatz, Galvin and Gagne ©2018

128



Dynamic Loading

- The entire program does need to be in memory to execute
- Routine is not loaded until it is called
- Better memory-space utilization; unused routine is never loaded
- All routines kept on disk in relocatable load format
- 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
 - OS can help by providing libraries to implement dynamic loading



Operating System Concepts – 10th Edition

9.129

Silberschatz, Galvin and Gagne ©2018

129



Dynamic Linking

- **Static linking** – system libraries and program code combined by the loader into the binary program image
- 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 checks if routine is in processes' memory address
 - If not in address space, add to address space
- Dynamic linking is particularly useful for libraries
- System also known as **shared libraries**
- Consider applicability to patching system libraries
 - Versioning may be needed

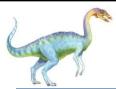


Operating System Concepts – 10th Edition

9.130

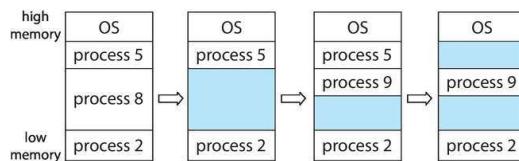
Silberschatz, Galvin and Gagne ©2018

130



Variable Partition

- Multiple-partition allocation
 - Degree of multiprogramming limited by number of partitions
 - **Variable-partition** sizes for efficiency (sized to a given process' needs)
 - **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
 - Process exiting frees its partition, adjacent free partitions combined
 - Operating system maintains information about:
 - a) allocated partitions
 - b) free partitions (hole)



Operating System Concepts – 10th Edition

9.131

Silberschatz, Galvin and Gagne ©2018

131

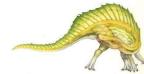


Dynamic Storage-Allocation Problem

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

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

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



Operating System Concepts – 10th Edition

9.132

Silberschatz, Galvin and Gagne ©2018

132



Fragmentation

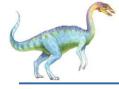
- **External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous
- **Internal Fragmentation** – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used
- First fit analysis reveals that given N blocks allocated, $0.5 N$ blocks lost to fragmentation
 - 1/3 may be unusable -> **50-percent rule**



Paging

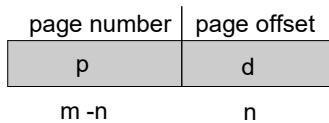
- Physical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available
 - Avoids external fragmentation
 - Avoids problem of varying sized memory chunks
- Divide physical memory into fixed-sized blocks called **frames**
 - Size is power of 2, between 512 bytes and 16 Mbytes
- 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
- Backing store likewise split into pages
- Still have Internal fragmentation





Address Translation Scheme

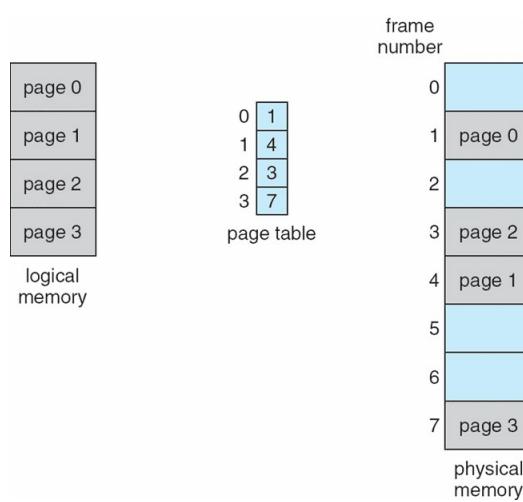
- Address generated by CPU is divided into:
 - **Page number (p)** – used as an index into a **page table** which contains base address of each page in physical memory
 - **Page offset (d)** – combined with base address to define the physical memory address that is sent to the memory unit

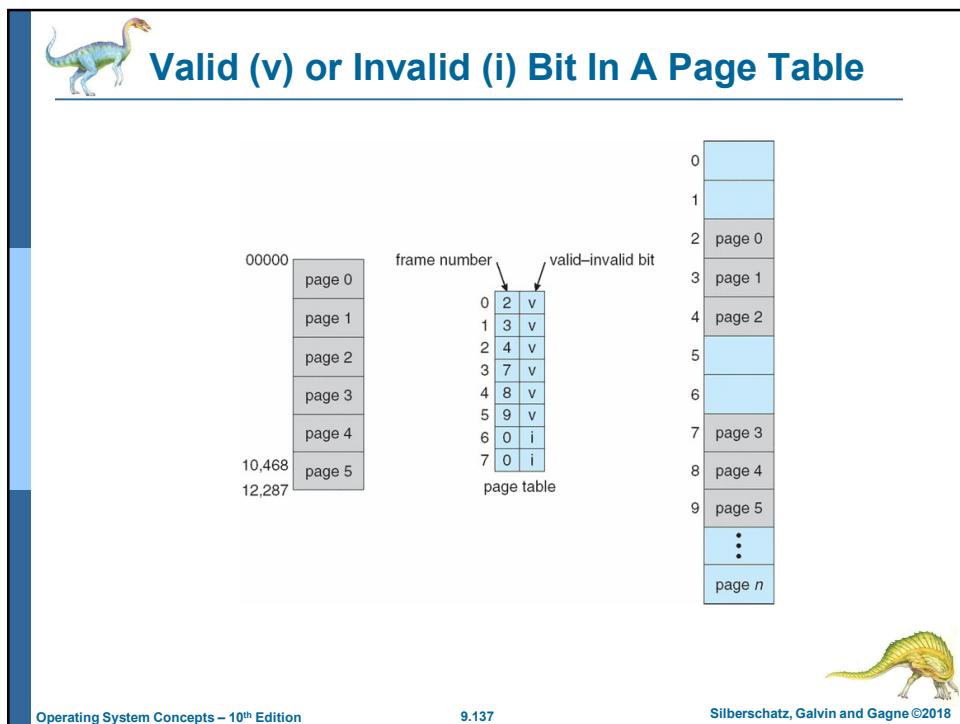


- For given logical address space 2^m and page size 2^n

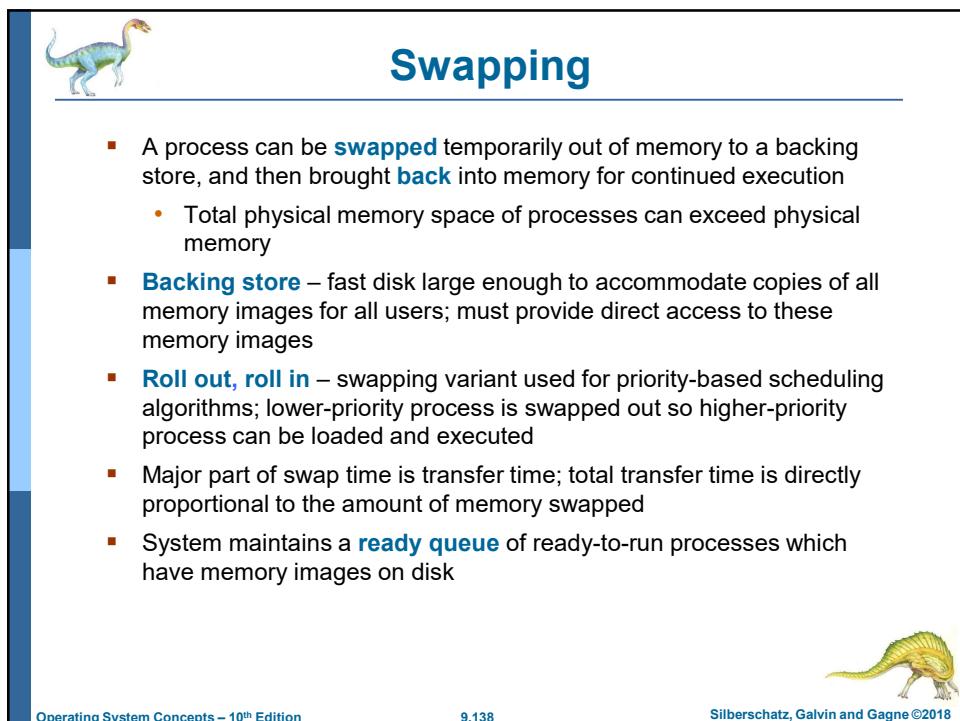


Paging Model of Logical and Physical Memory

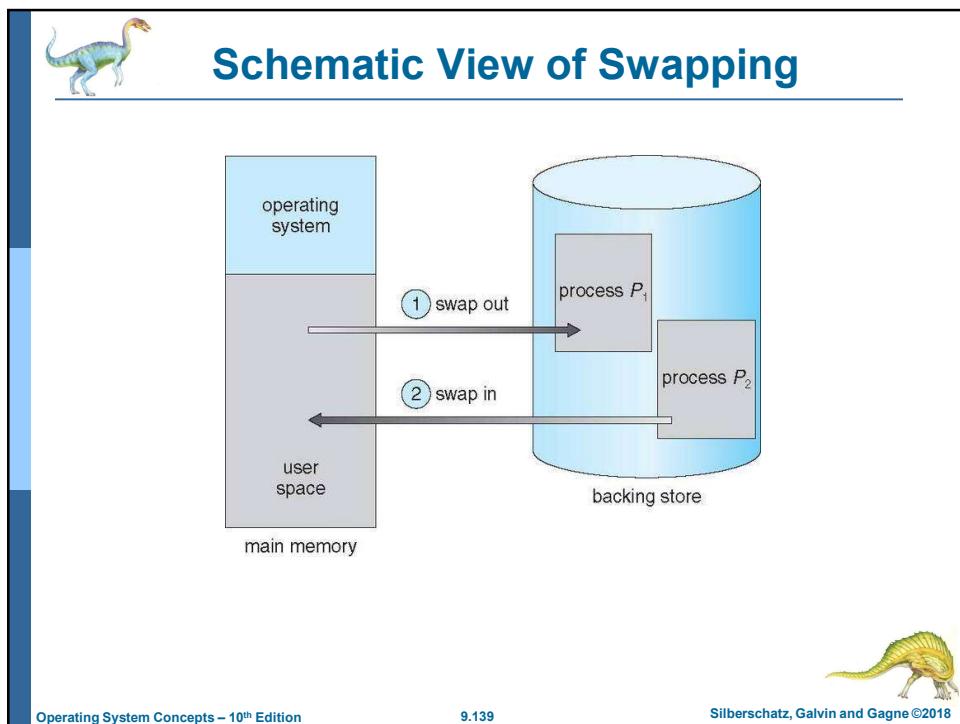




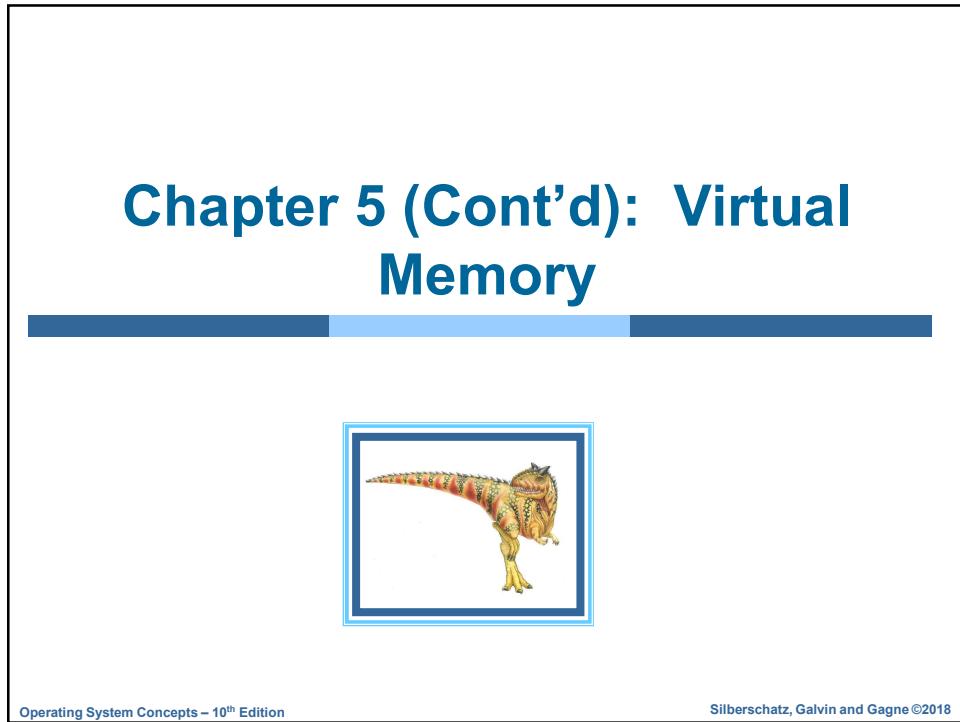
137



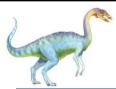
138



139



140



Virtual memory

- **Virtual memory** – separation of user logical memory from physical memory
 - Only part of the program needs to be in memory for execution
 - Logical address space can therefore be much larger than physical address space
 - Allows address spaces to be shared by several processes
 - Allows for more efficient process creation
 - More programs running concurrently
 - Less I/O needed to load or swap processes



Operating System Concepts – 10th Edition

9.141

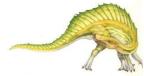
Silberschatz, Galvin and Gagne ©2018

141



Virtual memory (Cont.)

- **Virtual address space** – logical view of how process is stored in memory
 - Usually start at address 0, contiguous addresses until end of space
 - Meanwhile, physical memory organized in page frames
 - MMU must map logical to physical
- Virtual memory can be implemented via:
 - Demand paging
 - Demand segmentation



Operating System Concepts – 10th Edition

9.142

Silberschatz, Galvin and Gagne ©2018

142



Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed
 - Less I/O needed, no unnecessary I/O
 - Less memory needed
 - Faster response
 - More users
- Similar to paging system with swapping (diagram on right)
- Page is needed \Rightarrow reference to it
 - invalid reference \Rightarrow abort
 - not-in-memory \Rightarrow bring to memory
- **Lazy swapper** – never swaps a page into memory unless page will be needed
 - Swapper that deals with pages is a **pager**

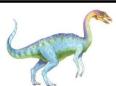


Operating System Concepts – 10th Edition

9.143

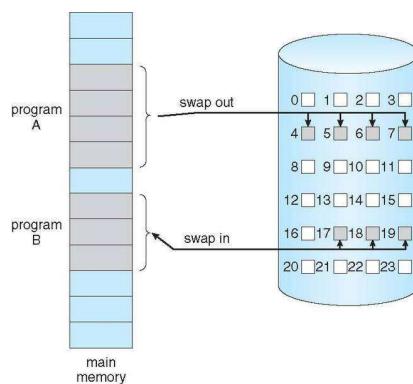
Silberschatz, Galvin and Gagne ©2018

143



Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed
 - Less I/O needed, no unnecessary I/O
 - Less memory needed
 - Faster response
 - More users
- Similar to paging system with swapping (diagram on right)



Operating System Concepts – 10th Edition

9.144

Silberschatz, Galvin and Gagne ©2018

144

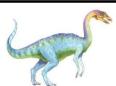


Valid-Invalid Bit

- With each page table entry a valid–invalid bit is associated (**v** ⇒ in-memory – **memory resident**, **i** ⇒ not-in-memory)
- Initially valid–invalid bit is set to **i** on all entries
- Example of a page table snapshot:

Frame #	valid-invalid bit
	v
	v
	v
	i
...	
	i
	i

- During MMU address translation, if valid–invalid bit in page table entry is **i** ⇒ page fault



Page Table When Some Pages Are Not in Main Memory

0	A
1	B
2	C
3	D
4	E
5	F
6	G
7	H

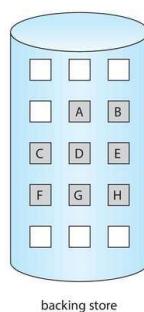
logical memory

frame	valid-invalid bit
0	4 v
1	i
2	6 v
3	i
4	i
5	9 v
6	i
7	i

page table

0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

physical memory

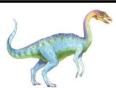


backing store



Copy-on-Write

- **Copy-on-Write** (COW) allows both parent and child processes to initially **share** the same pages in memory
 - If either process modifies a shared page, only then is the page copied
- COW allows more efficient process creation as only modified pages are copied
- In general, free pages are allocated from a **pool of zero-fill-on-demand** pages
 - Pool should always have free frames for fast demand page execution
 - Don't want to have to free a frame as well as other processing on page fault
 - Why zero-out a page before allocating it?
- `vfork()` variation on `fork()` system call has parent suspend and child using copy-on-write address space of parent
 - Designed to have child call `exec()`
 - Very efficient

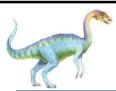


Basic Page Replacement

1. Find the location of the desired page on disk
2. Find a free frame:
 - If there is a free frame, use it
 - If there is no free frame, use a page replacement algorithm to select a **victim frame**
 - Write victim frame to disk if dirty
3. Bring the desired page into the (newly) free frame; update the page and frame tables
4. Continue the process by restarting the instruction that caused the trap

Note now potentially 2 page transfers for page fault – increasing EAT

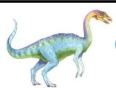




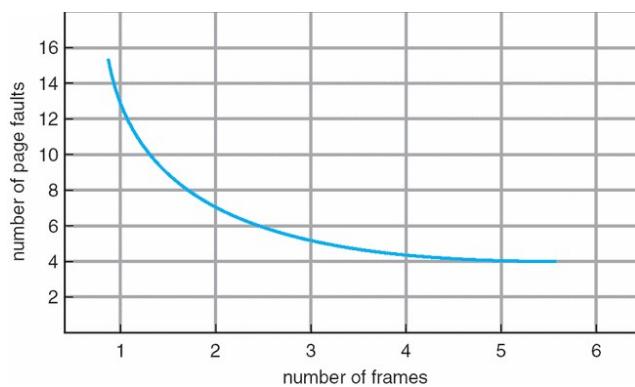
Page and Frame Replacement Algorithms

- **Frame-allocation algorithm** determines
 - How many frames to give each process
 - Which frames to replace
- **Page-replacement algorithm**
 - Want lowest page-fault rate on both first access and re-access
- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
 - String is just page numbers, not full addresses
 - Repeated access to the same page does not cause a page fault
 - Results depend on number of frames available
- In all our examples, the **reference string** of referenced page numbers is

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



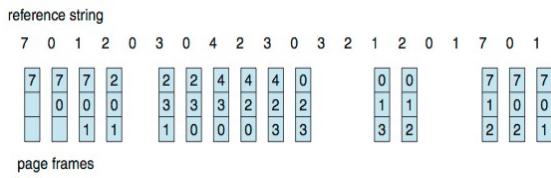
Graph of Page Faults Versus the Number of Frames





First-In-First-Out (FIFO) Algorithm

- Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**
- 3 frames (3 pages can be in memory at a time per process)

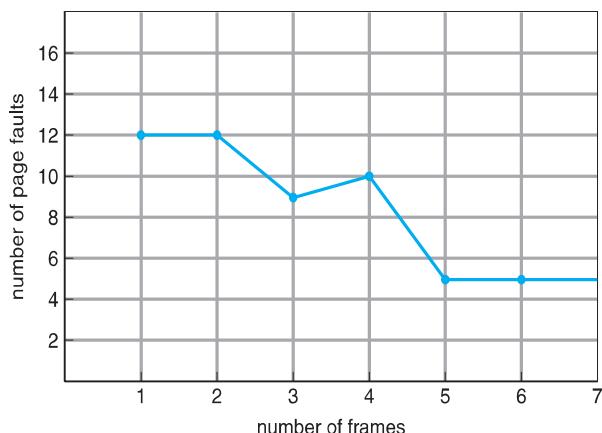


15 page faults

- Can vary by reference string: consider 1,2,3,4,1,2,5,1,2,3,4,5
 - Adding more frames can cause more page faults!
 - **Belady's Anomaly**
- How to track ages of pages?
 - Just use a FIFO queue



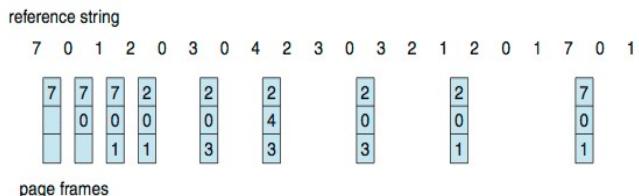
FIFO Illustrating Belady's Anomaly





Optimal Algorithm

- Replace page that will not be used for longest period of time
 - 9 is optimal for the example
- How do you know this?
 - Can't read the future
- Used for measuring how well your algorithm performs



Operating System Concepts – 10th Edition

9.153

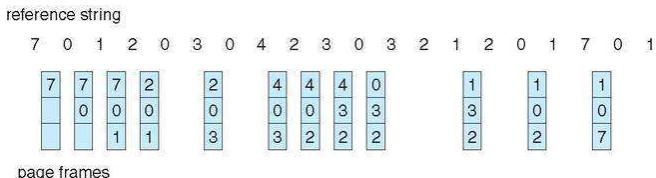
Silberschatz, Galvin and Gagne ©2018

153



Least Recently Used (LRU) Algorithm

- Use past knowledge rather than future
- Replace page that has not been used in the most amount of time
- Associate time of last use with each page

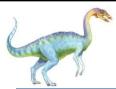


Operating System Concepts – 10th Edition

9.154

Silberschatz, Galvin and Gagne ©2018

154



LRU Algorithm (Cont.)

- 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 find smallest value
 - Search through table needed
- 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
 - But each update more expensive
 - No search for replacement

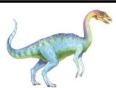


Operating System Concepts – 10th Edition

9.155

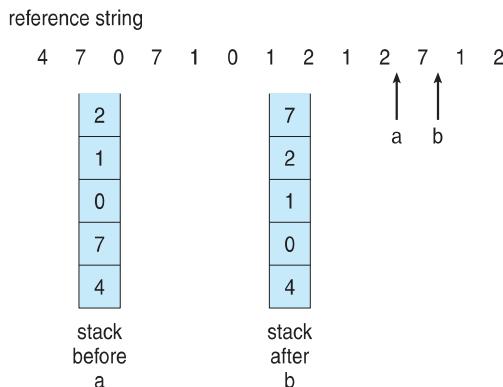
Silberschatz, Galvin and Gagne ©2018

155



LRU Algorithm (Cont.)

- LRU and OPT are cases of **stack algorithms** that don't have Belady's Anomaly
- Use Of A Stack to Record Most Recent Page References



Operating System Concepts – 10th Edition

9.156

Silberschatz, Galvin and Gagne ©2018

156



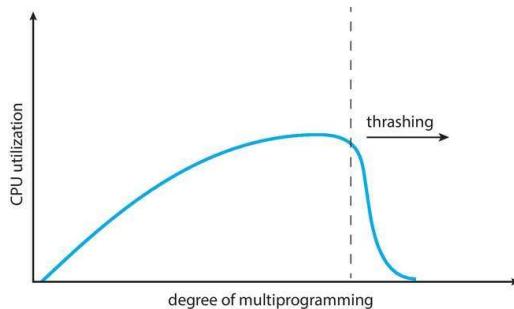
Thrashing

- If a process does not have “enough” pages, the page-fault rate is very high
 - Page fault to get page
 - Replace existing frame
 - But quickly need replaced frame back
 - This leads to:
 - Low CPU utilization
 - Operating system thinking that it needs to increase the degree of multiprogramming
 - Another process added to the system

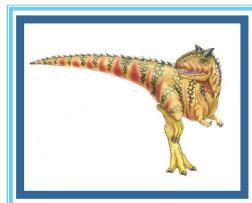


Thrashing (Cont.)

- **Thrashing.** A process is busy swapping pages in and out



Chapter 6: File Management



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

159

Outline

- File Concept
- Access Methods
- Disk and Directory Structure
- Protection
- Memory-Mapped Files



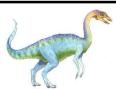
Operating System Concepts – 10th Edition

13.160

Silberschatz, Galvin and Gagne ©2018

160

80



File Concept

- Contiguous logical address space
- Types:
 - Data
 - Numeric
 - Character
 - Binary
 - Program
- Contents defined by file's creator
 - Many types
 - **text file**,
 - **source file**,
 - **executable file**



File Attributes

- **Name** – only information kept in human-readable form
- **Identifier** – unique tag (number) identifies file within file system
- **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
- Many variations, including extended file attributes such as file checksum
- Information kept in the directory structure





File Operations

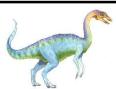
- **Create**
- **Write** – at **write pointer** location
- **Read** – at **read pointer** location
- **Reposition within file** - **seek**
- **Delete**
- **Truncate**
- **Open (F_i)** – search the directory structure on disk for entry F_i , and move the content of entry to memory
- **Close (F_i)** – move the content of entry F_i in memory to directory structure on disk



File Locking

- Provided by some operating systems and file systems
 - Similar to reader-writer locks
 - **Shared lock** similar to reader lock – several processes can acquire concurrently
 - **Exclusive lock** similar to writer lock
- Mediates access to a file
- Mandatory or advisory:
 - **Mandatory** – access is denied depending on locks held and requested
 - **Advisory** – processes can find status of locks and decide what to do





File Structure

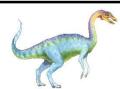
- None - sequence of words, bytes
- Simple record structure
 - Lines
 - Fixed length
 - Variable length
- Complex Structures
 - Formatted document
 - Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters
- Who decides:
 - Operating system
 - Program



Access Methods

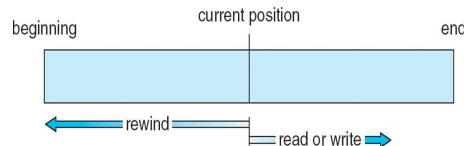
- A file is fixed length **logical records**
- **Sequential Access**
- **Direct Access**
- **Other Access Methods**





Sequential Access

- Operations
 - **read next**
 - **write next**
 - **Reset**
 - no read after last write (rewrite)
- Figure



Operating System Concepts – 10th Edition

13.167

Silberschatz, Galvin and Gagne ©2018



167



Direct Access

- Operations
 - **read *n***
 - **write *n***
 - **position to *n***
 - **read next**
 - **write next**
 - **rewrite *n***
- $n = \text{relative block number}$
- Relative block numbers allow OS to decide where file should be placed

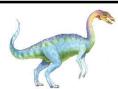
Operating System Concepts – 10th Edition

13.168

Silberschatz, Galvin and Gagne ©2018

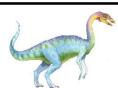


168

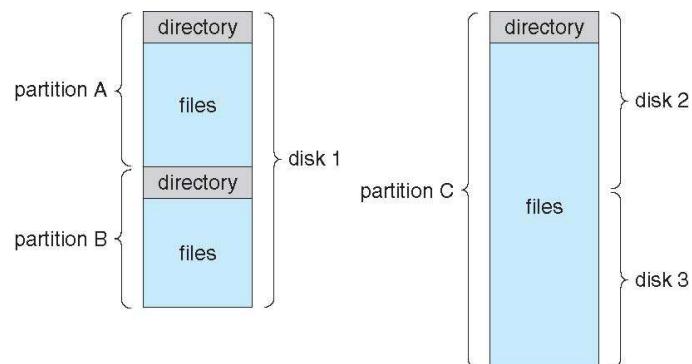


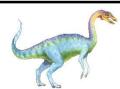
Disk Structure

- Disk can be subdivided into **partitions**
- Disks or partitions can be **RAID** protected against failure
- Disk or partition can be used **raw** – without a file system, or **formatted** with a file system
- Partitions also known as minidisks, slices
- Entity containing file system is known as a **volume**
- Each volume containing a file system also tracks that file system's info in **device directory** or **volume table of contents**
- In addition to **general-purpose file systems** there are many **special-purpose file systems**, frequently all within the same operating system or computer



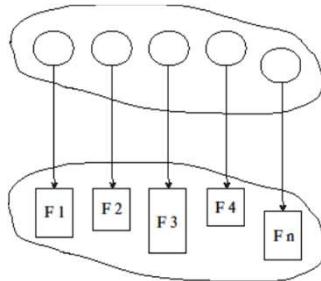
A Typical File-system Organization





Directory Structure

- A collection of nodes containing information about all files



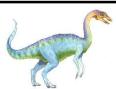
- Both the directory structure and the files reside on disk



Operations Performed on Directory

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





Directory Organization

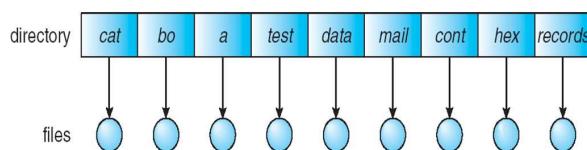
The directory is organized logically to obtain

- 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 Java programs, all games, ...)



Single-Level Directory

- A single directory for all users



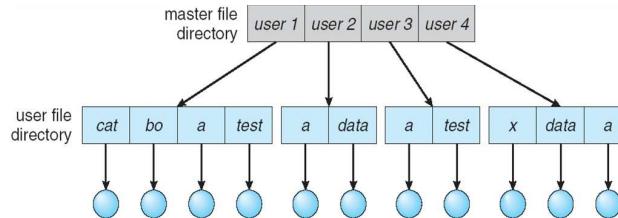
- Naming problem
- Grouping problem





Two-Level Directory

- Separate directory for each user

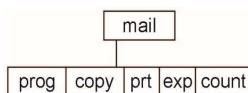


- Path name
- Can have the same file name for different user
- Efficient searching
- No grouping capability



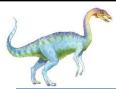
Current Directory

- Can designate one of the directories as the current (working) directory
 - `cd /spell/mail/prog`
 - `type list`
- Creating and deleting a file is done in current directory
- Example of creating a new file
 - If in current directory is `/mail`
 - The command
`mkdir <dir-name>`
 - Results in:



- Deleting "mail" ⇒ deleting the entire subtree rooted by "mail"



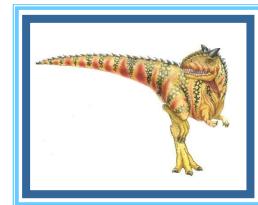


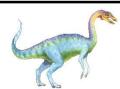
Protection

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



Chapter 6 (Cont'd): File System Implementation



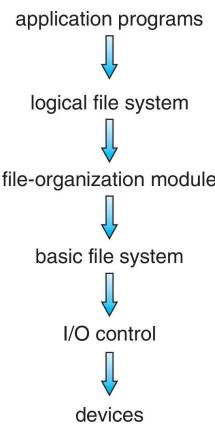


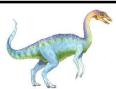
File-System Structure

- File structure
 - Logical storage unit
 - Collection of related information
- **File system** resides on secondary storage (disks)
 - Provided user interface to storage, mapping logical to physical
 - Provides efficient and convenient access to disk by allowing data to be stored, located retrieved easily
- Disk provides in-place rewrite and random access
 - I/O transfers performed in **blocks of sectors** (usually 512 bytes)
- **File control block (FCB)** – storage structure consisting of information about a file
- **Device driver** controls the physical device
- File system organized into layers



Layered File System





File System Layers

- **Device drivers** manage I/O devices at the I/O control layer
 - Given commands like
 - read drive1, cylinder 72, track 2, sector 10, into memory location 1060
 - Outputs low-level hardware specific commands to hardware controller
- **Basic file system** given command like “retrieve block 123” translates to device driver
 - Also manages memory buffers and caches (allocation, freeing, replacement)
 - Buffers hold data in transit
 - Caches hold frequently used data
 - **File organization module** understands files, logical address, and physical blocks
 - Translates logical block # to physical block #
 - Manages free space, disk allocation



File System Layers (Cont.)

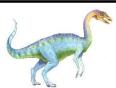
- **Logical file system** manages metadata information
 - Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)
 - Directory management
 - Protection
 - Layering useful for reducing complexity and redundancy, but adds overhead and can decrease performance
 - Logical layers can be implemented by any coding method according to OS designer





Allocation Method

- An allocation method refers to how disk blocks are allocated for files:
 - Contiguous
 - Linked
 - File Allocation Table (FAT)



Contiguous Allocation Method

- An allocation method refers to how disk blocks are allocated for files:
- Each file occupies set of contiguous blocks
 - Best performance in most cases
 - Simple – only starting location (block #) and length (number of blocks) are required
 - Problems include:
 - ▶ Finding space on the disk for a file,
 - ▶ Knowing file size,
 - ▶ External fragmentation, need for **compaction off-line** (**downtime**) or **on-line**

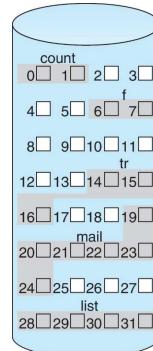




Contiguous Allocation (Cont.)

- Mapping from logical to physical
(block size = 512 bytes)

LA/512
Q
R



directory		
file	start	length
count	0	2
0	1	3
f	4	5
tr	6	7
mail	8	9
list	10	11
tr	12	13
mail	14	15
list	16	17
list	18	19
mail	20	21
list	22	23
list	24	25
list	26	27
list	28	29
list	30	31



Operating System Concepts – 10th Edition

13.185

Silberschatz, Galvin and Gagne ©2018

185



Linked Allocation

- Each file is a linked list of blocks
- File ends at nil pointer
- No external fragmentation
- Each block contains pointer to next block
- No compaction, external fragmentation
- Free space management system called when new block needed
- Improve efficiency by clustering blocks into groups but increases internal fragmentation
- Reliability can be a problem
- Locating a block can take many I/Os and disk seeks



Operating System Concepts – 10th Edition

13.186

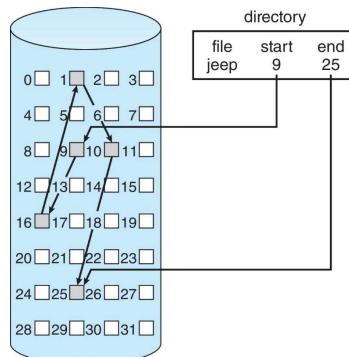
Silberschatz, Galvin and Gagne ©2018

186



Linked Allocation Example

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



Operating System Concepts – 10th Edition

13.187

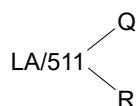
Silberschatz, Galvin and Gagne ©2018

187



Linked Allocation (Cont.)

- Mapping



- Block to be accessed is the Q^{th} block in the linked chain of blocks representing the file.
- Displacement into block = $R + 1$



Operating System Concepts – 10th Edition

13.188

Silberschatz, Galvin and Gagne ©2018

188

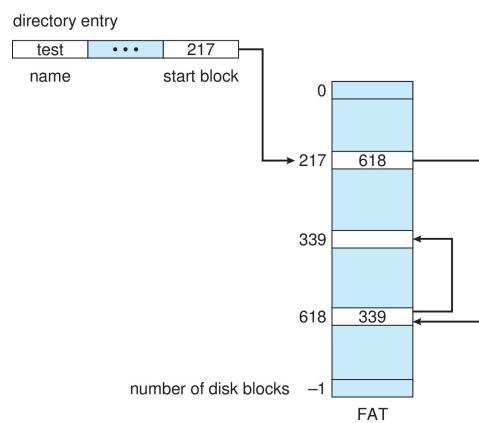


FAT Allocation Method

- Beginning of volume has table, indexed by block number
- Much like a linked list, but faster on disk and cacheable
- New block allocation simple



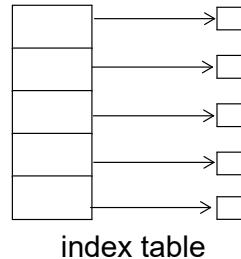
File-Allocation Table





Indexed Allocation Method

- Each file has its own **index block**(s) of pointers to its data blocks
- Logical view



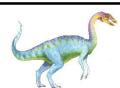
Operating System Concepts – 10th Edition

13.191

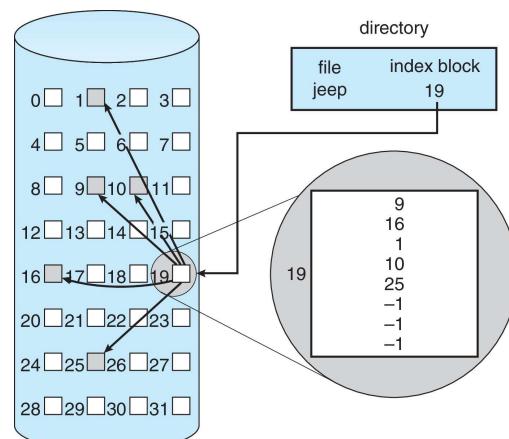
Silberschatz, Galvin and Gagne ©2018



191



Example of Indexed Allocation



Operating System Concepts – 10th Edition

13.192

Silberschatz, Galvin and Gagne ©2018



192



File Sharing

- Allows multiple users / systems access to the same files
- Permissions / protection must be implemented and accurate
 - Most systems provide concepts of owner, group member
 - Must have a way to apply these between systems



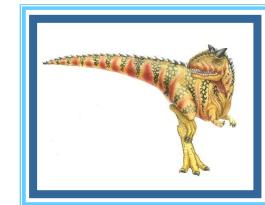
Operating System Concepts – 10th Edition

13.193

Silberschatz, Galvin and Gagne ©2018

193

Chapter 6 (Cont'd): OS Security



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018

194



The Security Problem

- System **secure** if resources used and accessed as intended under all circumstances
 - Unachievable
- **Intruders (crackers)** attempt to breach security
- **Threat** is potential security violation
- **Attack** is attempt to breach security
- Attack can be accidental or malicious
- Easier to protect against accidental than malicious misuse



Operating System Concepts – 10th Edition

13.195

Silberschatz, Galvin and Gagne ©2018

195



Security Violation Categories

- **Breach of confidentiality**
 - Unauthorized reading of data
- **Breach of integrity**
 - Unauthorized modification of data
- **Breach of availability**
 - Unauthorized destruction of data
- **Theft of service**
 - Unauthorized use of resources
- **Denial of service (DOS)**
 - Prevention of legitimate use

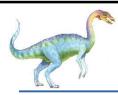


Operating System Concepts – 10th Edition

13.196

Silberschatz, Galvin and Gagne ©2018

196



Security Violation Methods

- **Masquerading** (breach **authentication**)
 - Pretending to be an authorized user to escalate privileges
- **Replay attack**
 - As is or with **message modification**
- **Man-in-the-middle attack**
 - Intruder sits in data flow, masquerading as sender to receiver and vice versa
- **Session hijacking**
 - Intercept an already-established session to bypass authentication
- **Privilege escalation**
 - Common attack type with access beyond what a user or resource is supposed to have



Operating System Concepts – 10th Edition

13.197

Silberschatz, Galvin and Gagne ©2018

197



Security Measure Levels

- Impossible to have absolute security, but make cost to perpetrator sufficiently high to deter most intruders
- Security must occur at four levels to be effective:
 - **Physical**
 - Data centers, servers, connected terminals
 - **Application**
 - Benign or malicious apps can cause security problems
 - **Operating System**
 - Protection mechanisms, debugging
 - **Network**
 - Intercepted communications, interruption, DOS
- Security is as weak as the weakest link in the chain
- Humans a risk too via **phishing** and **social-engineering** attacks
- But can too much security be a problem?



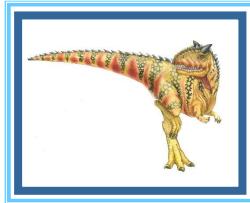
Operating System Concepts – 10th Edition

13.198

Silberschatz, Galvin and Gagne ©2018

198

End of Chapter 6



Operating System Concepts – 10th Edition

Silberschatz, Galvin and Gagne ©2018