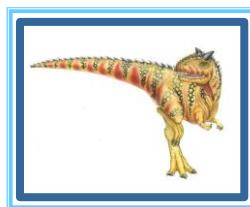


# Chapter 1: Introduction



Operating System Concepts – 10<sup>th</sup> 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 – 10<sup>th</sup> Edition

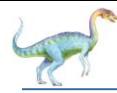
1.2

Silberschatz, Galvin and Gagne ©2018



2

1



## What is an Operating System?

### definition ↗

- 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



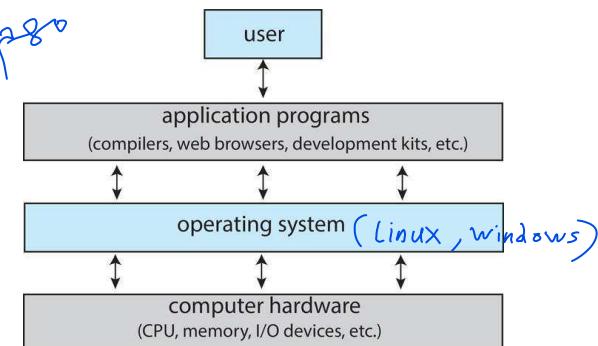
## Computer System Structure

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





## Abstract View of Components of Computer



Operating System Concepts – 10<sup>th</sup> Edition

1.5

Silberschatz, Galvin and Gagne ©2018



## What Operating Systems Do

- Depends on the point of view
  - Users want convenience, ease of use and good performance
- Users* Don't care about resource utilization
- But shared computer such as mainframe or minicomputer must keep all users happy
    - Operating system is a resource allocator and control program making efficient use of HW and managing execution of user programs
  - Users of dedicated systems such as workstations have dedicated resources but frequently use shared resources from servers
  - Mobile devices like smartphones and tablets are resource poor, optimized for usability and battery life
    - Mobile user interfaces such as touch screens, voice recognition
  - Some computers have little or no user interface, such as embedded computers in devices and automobiles
    - Run primarily without user intervention

Operating System Concepts – 10<sup>th</sup> Edition

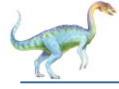
1.6

Silberschatz, Galvin and Gagne ©2018



6

3



## Overview of Computer System Structure



Operating System Concepts – 10<sup>th</sup> Edition

1.7

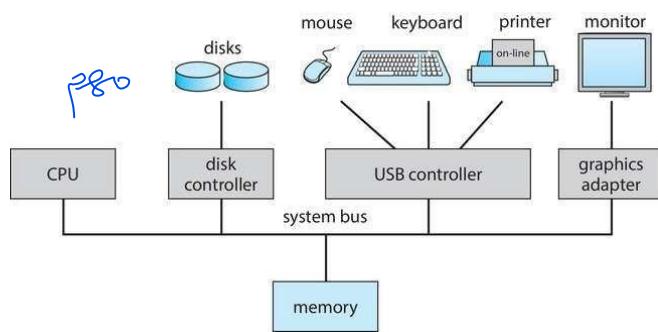
Silberschatz, Galvin and Gagne ©2018

7



## Computer System Organization

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



Operating System Concepts – 10<sup>th</sup> 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 Concepts – 10<sup>th</sup> Edition

1.9

Silberschatz, Galvin and Gagne ©2018

9



## Operating-System Operations

- Bootstrap program – simple code to initialize the system, load the kernel
- Kernel loads
- **I** Starts **system daemons** (services provided outside of the kernel)
- **2** **Kernel interrupt driven** (hardware and software)
  - i. Hardware interrupt by one of the devices
  - ii. 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



Operating System Concepts – 10<sup>th</sup> Edition

1.10

Silberschatz, Galvin and Gagne ©2018

10



## 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



Operating System Concepts – 10<sup>th</sup> Edition

1.11

Silberschatz, Galvin and Gagne ©2018

11



## Multitasking (Timesharing)

*definition :-*

- 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



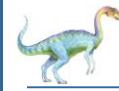
Operating System Concepts – 10<sup>th</sup> Edition

1.12

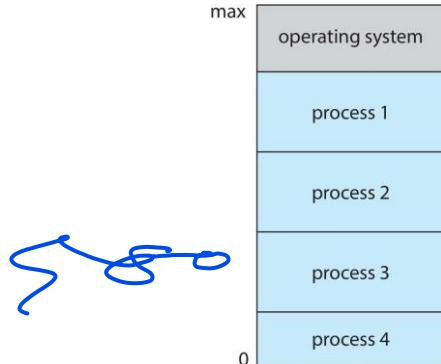
Silberschatz, Galvin and Gagne ©2018

12

~91>



## Memory Layout for Multiprogrammed System



Operating System Concepts – 10<sup>th</sup> Edition

1.13

Silberschatz, Galvin and Gagne ©2018

13



## Dual-mode Operation

- Dual-mode operation allows OS to protect itself and other system components
  - User mode and kernel mode
- Mode bit provided by hardware
  - Provides ability to distinguish when system is running user code or kernel code.
  - When a user is running ⇒ mode bit is "user"
  - When kernel code is executing ⇒ mode bit is "kernel"
- How do we guarantee that user does not explicitly set the mode bit to "kernel"?
  - System call changes mode to kernel, return from call resets it to user
- Some instructions designated as privileged, only executable in kernel mode



Operating System Concepts – 10<sup>th</sup> Edition

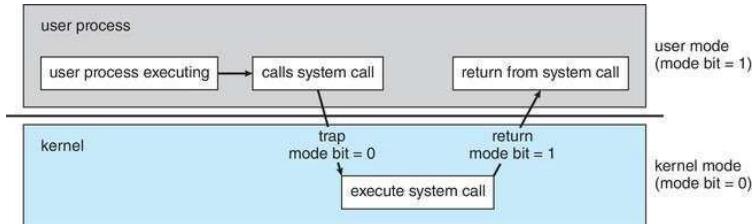
1.14

Silberschatz, Galvin and Gagne ©2018

14



## Transition from User to Kernel Mode



280



Operating System Concepts – 10<sup>th</sup> Edition

1.15

Silberschatz, Galvin and Gagne ©2018

15



## Virtualization

### definition:-

- Allows operating systems to run applications within other OSes
  - Vast and growing industry
- Emulation used when source CPU type different from target type (i.e. PowerPC to Intel x86)
  - Generally slowest method
  - When computer language not compiled to native code – Interpretation
- Virtualization – OS natively compiled for CPU, running guest OSes also natively compiled
  - Consider VMware running WinXP guests, each running applications, all on native WinXP host OS
  - VMM (virtual machine Manager) provides virtualization services

like VMware



Operating System Concepts – 10<sup>th</sup> Edition

1.16

Silberschatz, Galvin and Gagne ©2018

16



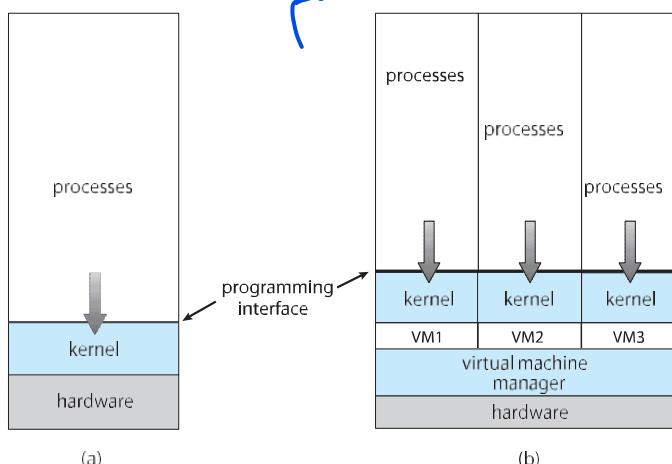
## Virtualization (cont.)

- **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

280





(Self Study)

## Computer System Architecture



Operating System Concepts – 10<sup>th</sup> Edition

1.19

Silberschatz, Galvin and Gagne ©2018

19



## Computer System Architecture (Self Study)

- Most systems use a single general-purpose processor
  - Most systems have special-purpose processors as well
- **Multiprocessors** systems growing in use and importance
  - Also known as **parallel systems, tightly-coupled systems**
  - Advantages include:
    1. **Increased throughput**
    2. **Economy of scale**
    3. **Increased reliability** – graceful degradation or fault tolerance
  - Two types:
    1. **Asymmetric Multiprocessing** – each processor is assigned a specific task.
    2. **Symmetric Multiprocessing** – each processor performs all tasks



Operating System Concepts – 10<sup>th</sup> Edition

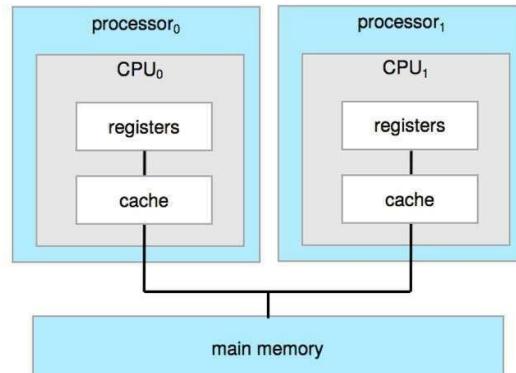
1.20

Silberschatz, Galvin and Gagne ©2018

20

10

## Symmetric Multiprocessing Architecture (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

1.21

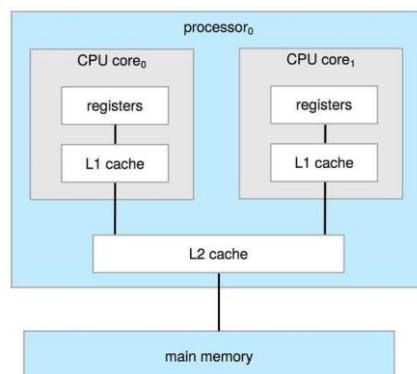
Silberschatz, Galvin and Gagne ©2018

21

## Dual-Core Design (Self Study)



- Multi-chip and **multicore**
- Systems containing all chips
  - Chassis containing multiple separate systems



Operating System Concepts – 10<sup>th</sup> Edition

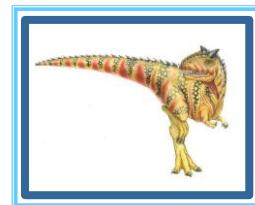
1.22

Silberschatz, Galvin and Gagne ©2018

22

11

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



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

23



## Outline

- Operating System Services
- User and Operating System-Interface
- System Calls
- Operating System Structure

Operating System Concepts – 10<sup>th</sup> Edition

1.24

Silberschatz, Galvin and Gagne ©2018



24

12



## 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:
  - **User interface** - Almost all operating systems have a user interface ([UI](#)).
    - Varies between [Command-Line \(CLI\)](#), [Graphics User Interface \(GUI\)](#), [touch-screen](#), [Batch](#)
  - **Program execution** - The system must be able to load a program into memory and to run that program, end execution, either normally or abnormally (indicating error)
  - **I/O operations** - A running program may require I/O, which may involve a file or an I/O device
  - **File-system manipulation** - The file system is of particular interest. Programs need to read and write files and directories, create and delete them, search them, list file information, permission management.



## Operating System Services (Cont.)

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





## Operating System Services (Cont.)

- Another set of OS functions exists for ensuring the efficient operation of the system itself via resource sharing
  - **Resource allocation** - When multiple users or multiple jobs running concurrently, resources must be allocated to each of them
    - Many types of resources - CPU cycles, main memory, file storage, I/O devices.
  - **Logging** - To keep track of which users use how much and what kinds of computer resources
  - **Protection and security** - The owners of information stored in a multiuser or networked computer system may want to control use of that information, concurrent processes should not interfere with each other
    - **Protection** involves ensuring that all access to system resources is controlled
    - **Security** of the system from outsiders requires user authentication, extends to defending external I/O devices from invalid access attempts



Operating System Concepts – 10<sup>th</sup> Edition

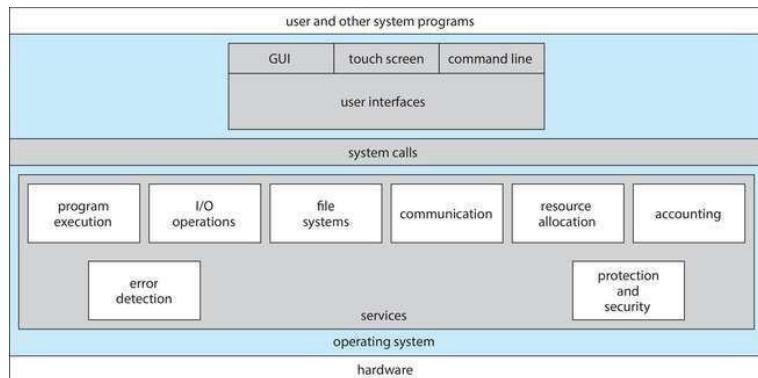
1.27

Silberschatz, Galvin and Gagne ©2018

27



## A View of Operating System Services



Operating System Concepts – 10<sup>th</sup> Edition

1.28

Silberschatz, Galvin and Gagne ©2018

28

14

## Command Line interpreter (Self Study)

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



Operating System Concepts – 10<sup>th</sup> Edition

1.29

Silberschatz, Galvin and Gagne ©2018

29

## Bourne Shell Command Interpreter (Self Study)



```
1. root@r6181-d5-us01:~ (ssh)
root@r6181-d5-us01:~ %11 ssh %22 %32 root@r6181-d5-us01... %3
Last login: Thu Jul 14 08:47:01 on ttys002
iMacPro:~ pbg$ ssh root@r6181-d5-us01
root@r6181-d5-us01's password:
Last login: Thu Jul 14 06:01:11 2016 from 172.16.16.162
[root@r6181-d5-us01 ~]# uptime
 06:57:48 up 16 days, 10:52,  3 users,  load average: 129.52, 80.33, 56.55
[root@r6181-d5-us01 ~]# df -kh
Filesystem      Size  Used Avail Use% Mounted on
/dev/mapper/vg_ks-lv_root
                  50G   19G   28G  41% /
tmpfs           127G  520K  127G   1% /dev/shm
/dev/sda1        477M   71M   381M  16% /boot
/dev/dssd0000    1.0T  480G  545G  47% /dssd_xfs
tcp://192.168.150.1:3334/orangefs
                  12T  5.7T  6.4T  47% /mnt/orangefs
/dev/gpfstest    23T  1.1T   22T   5% /mnt/gpfst
[root@r6181-d5-us01 ~]#
[root@r6181-d5-us01 ~]# ps aux | sort -nrk 3,3 | head -n 5
root     97653 11.2  6.6 42665344 17520636 ?  S<Ll Jul13 166:23 /usr/lpp/mmfs/bin/mmfsd
root     69849  6.6  0.0     0  0 ?  S  Jul12 181:54 [vpthread-1-1]
root     69850  6.4  0.0     0  0 ?  S  Jul12 177:42 [vpthread-1-2]
root     3829  3.0  0.0     0  0 ?  S  Jun27 730:04 [rp_thread 7:0]
root     3826  3.0  0.0     0  0 ?  S  Jun27 728:08 [rp_thread 6:0]
[root@r6181-d5-us01 ~]# ls -l /usr/lpp/mmfs/bin/mmfsd
-r-x----- 1 root root 20667161 Jun  3  2015 /usr/lpp/mmfs/bin/mmfsd
[root@r6181-d5-us01 ~]#
```

Operating System Concepts – 10<sup>th</sup> Edition

1.30

Silberschatz, Galvin and Gagne ©2018

30

15

## User Operating System Interface - GUI (Self Study)

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



Operating System Concepts – 10<sup>th</sup> Edition

1.31

Silberschatz, Galvin and Gagne ©2018

31

## Touchscreen Interfaces (Self Study)

- Touchscreen devices require new interfaces
  - Mouse not possible or not desired
  - Actions and selection based on gestures
  - Virtual keyboard for text entry
- Voice commands



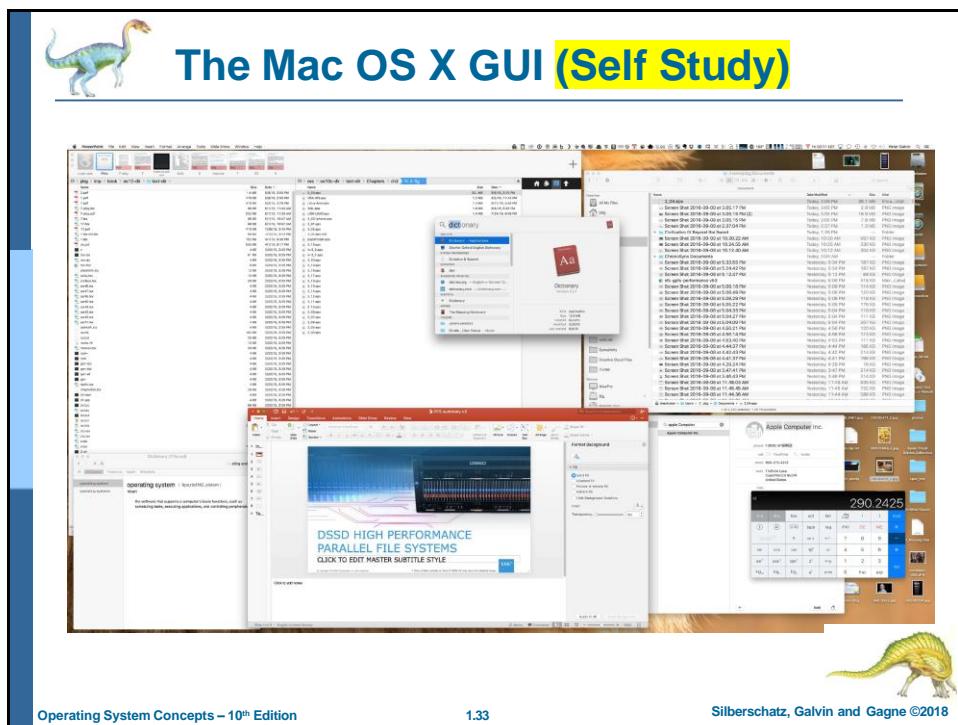
Operating System Concepts – 10<sup>th</sup> Edition

1.32

Silberschatz, Galvin and Gagne ©2018

32

16



33

## System Calls (Self Study)

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

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

Operating System Concepts – 10<sup>th</sup> Edition

1.34

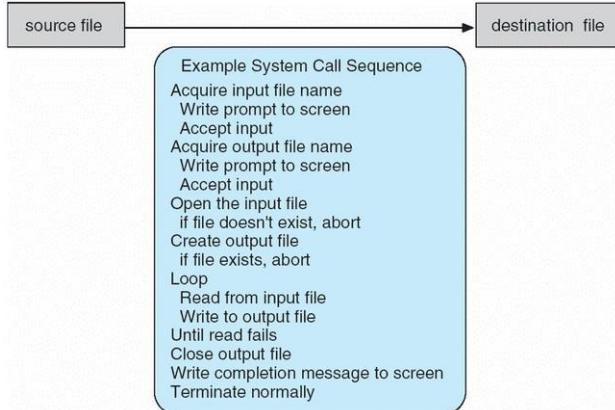
Silberschatz, Galvin and Gagne ©2018

34



## Example of System Calls (Self Study)

- System call sequence to copy the contents of one file to another file



Operating System Concepts – 10<sup>th</sup> Edition

1.35

Silberschatz, Galvin and Gagne ©2018



35



## System Call Implementation (Self Study)

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

Operating System Concepts – 10<sup>th</sup> Edition

1.36

Silberschatz, Galvin and Gagne ©2018



36



## Types of System Calls (Self Study)

- Process control
  - create process, terminate process
  - end, abort
  - load, execute
  - get process attributes, set process attributes
  - wait for time
  - wait event, signal event
  - allocate and free memory
  - Dump memory if error
  - **Debugger** for determining **bugs**, **single step** execution
  - **Locks** for managing access to shared data between processes



## Types of System Calls (Cont.) (Self Study)

- File management
  - create file, delete file
  - open, close file
  - read, write, reposition
  - get and set file attributes
- Device management
  - request device, release device
  - read, write, reposition
  - get device attributes, set device attributes
  - logically attach or detach devices





## Types of System Calls (Cont.) (Self Study)

- Information maintenance
  - get time or date, set time or date
  - get system data, set system data
  - get and set process, file, or device attributes
- Communications
  - create, delete communication connection
  - send, receive messages if **message passing model** to **host name** or **process name**
    - From **client** to **server**
  - **Shared-memory model** create and gain access to memory regions
  - transfer status information
  - attach and detach remote devices



Operating System Concepts – 10<sup>th</sup> Edition

1.39

Silberschatz, Galvin and Gagne ©2018

39



## Types of System Calls (Cont.) (Self Study)

- Protection
  - Control access to resources
  - Get and set permissions
  - Allow and deny user access



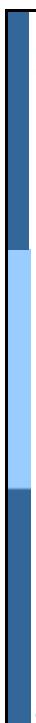
Operating System Concepts – 10<sup>th</sup> Edition

1.40

Silberschatz, Galvin and Gagne ©2018

40

20

## Examples of Windows and Unix System Calls (Self Study)

**EXAMPLES OF WINDOWS AND UNIX SYSTEM CALLS**

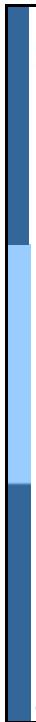
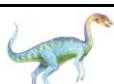
The following illustrates various equivalent system calls for Windows and UNIX operating systems.

|                         | Windows   | Unix                                   |
|-------------------------|---|--|
| Process control         | CreateProcess()<br>ExitProcess()<br>WaitForSingleObject()                           | fork()<br>exit()<br>wait()             |
| File management         | CreateFile()<br>ReadFile()<br>WriteFile()<br>CloseHandle()                          | open()<br>read()<br>write()<br>close() |
| Device management       | SetConsoleMode()<br>ReadConsole()<br>WriteConsole()                                 | ioctl()<br>read()<br>write()           |
| Information maintenance | GetCurrentProcessID()<br>SetTimer()<br>Sleep()                                      | getpid()<br>alarm()<br>sleep()         |
| Communications          | CreatePipe()<br>CreatefileMapping()<br>MapViewOfFile()                              | pipe()<br>shm_open()<br>mmap()         |
| Protection              | SetFileSecurity()<br>InitializeSecurityDescriptor()<br>SetSecurityDescriptorGroup() | chmod()<br>umask()<br>chown()          |



**Operating System Concepts – 10<sup>th</sup> Edition**      **1.41**      **Silberschatz, Galvin and Gagne ©2018**

41

## Operating System Structure

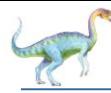
- General-purpose OS is very large program
- Various ways to structure ones
  - Simple structure – MS-DOS
  - More complex – UNIX
  - Layered – an abstraction
  - Microkernel – Mach



**Operating System Concepts – 10<sup>th</sup> Edition**      **1.42**      **Silberschatz, Galvin and Gagne ©2018**

42

12/04



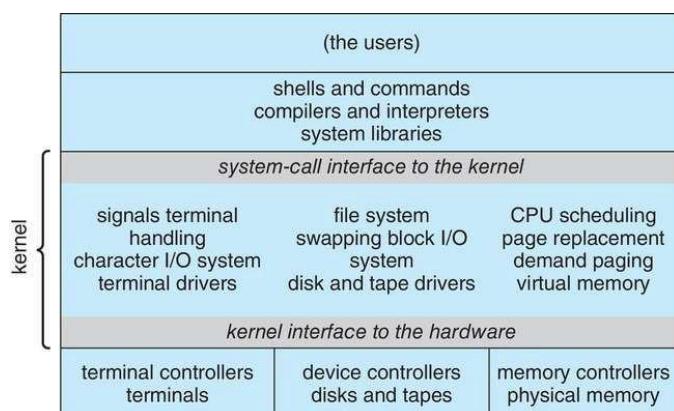
## Monolithic Structure – Original UNIX

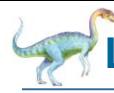
- UNIX – limited by hardware functionality, the original UNIX operating system had limited structuring.
- The UNIX OS consists of two separable parts
  - Systems programs
  - The kernel
    - Consists of everything below the system-call interface and above the physical hardware
    - Provides the file system, CPU scheduling, memory management, and other operating-system functions; a large number of functions for one level



## Traditional UNIX System Structure (Self Study)

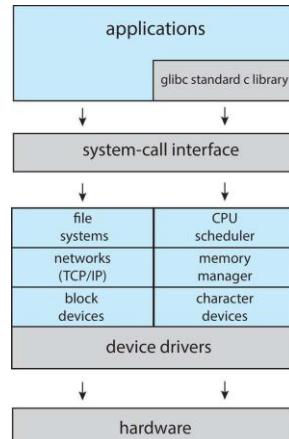
Beyond simple but not fully layered





## Linux System Structure (Self Study)

Monolithic plus modular design



Operating System Concepts – 10<sup>th</sup> Edition

1.45

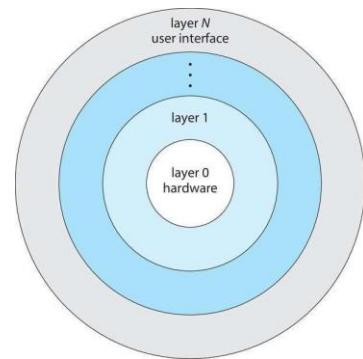
Silberschatz, Galvin and Gagne ©2018

45



## 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 – 10<sup>th</sup> Edition

1.46

Silberschatz, Galvin and Gagne ©2018

46



## 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 – 10<sup>th</sup> Edition

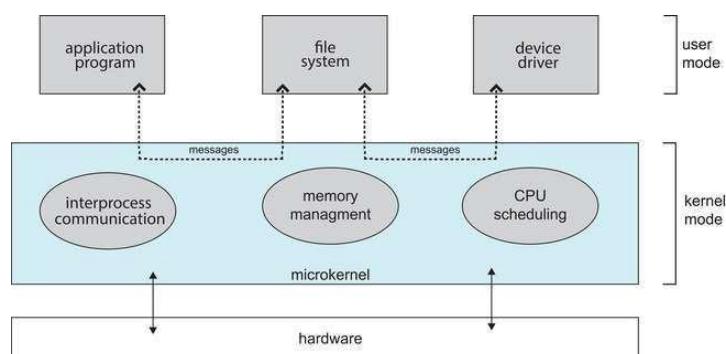
1.47

Silberschatz, Galvin and Gagne ©2018

47



## Microkernel System Structure

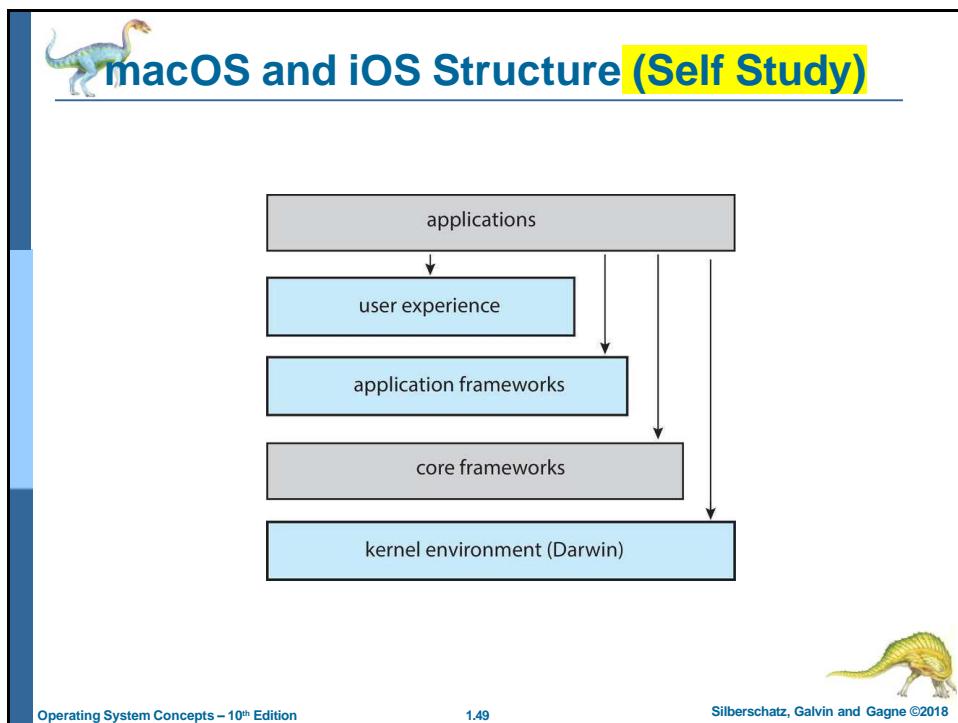


Operating System Concepts – 10<sup>th</sup> Edition

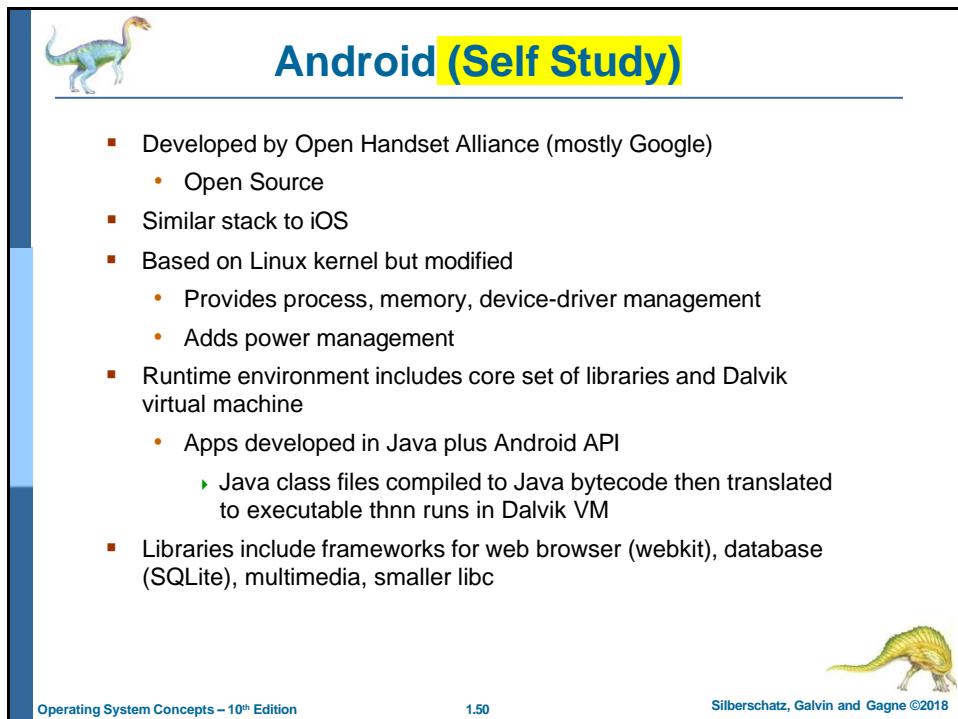
1.48

Silberschatz, Galvin and Gagne ©2018

48



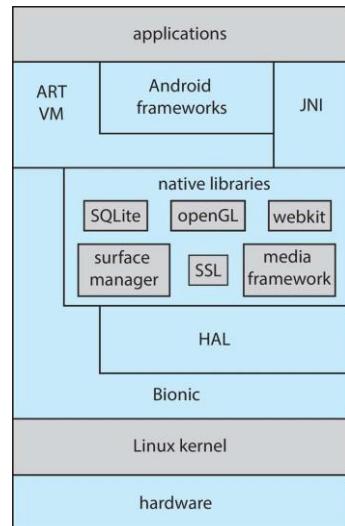
49



50



## Android Architecture (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

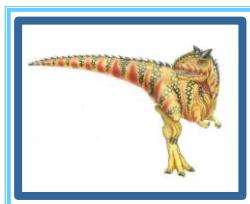
1.51

Silberschatz, Galvin and Gagne ©2018



51

## End of Chapter 1



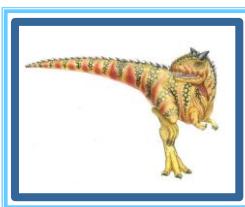
Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

52

26

# Chapter 2: Processes & Threads



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

53



## Outline

- Process Concept
- Process Scheduling
- Operations on Processes

Operating System Concepts – 10<sup>th</sup> Edition

3.54

Silberschatz, Galvin and Gagne ©2018

54





## Process Concept

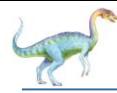
- An operating system executes a variety of programs that run as a process.
- Process – a program in execution; process execution must progress in sequential fashion. No parallel execution of instructions of a single process
- Multiple parts *of Process*
  1. The program code, also called **text section**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



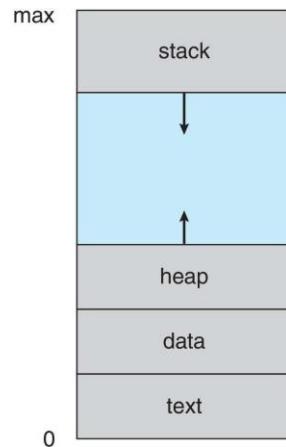
## Process Concept (Cont.)

- 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 in Memory (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

3.57

Silberschatz, Galvin and Gagne ©2018

57



## 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**



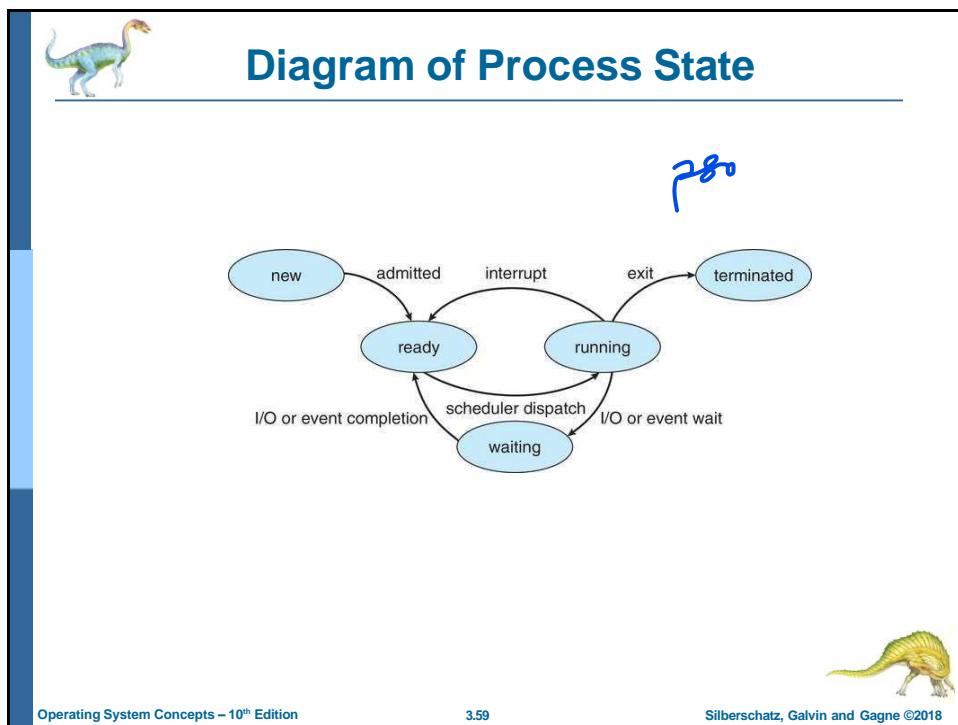
Operating System Concepts – 10<sup>th</sup> Edition

3.58

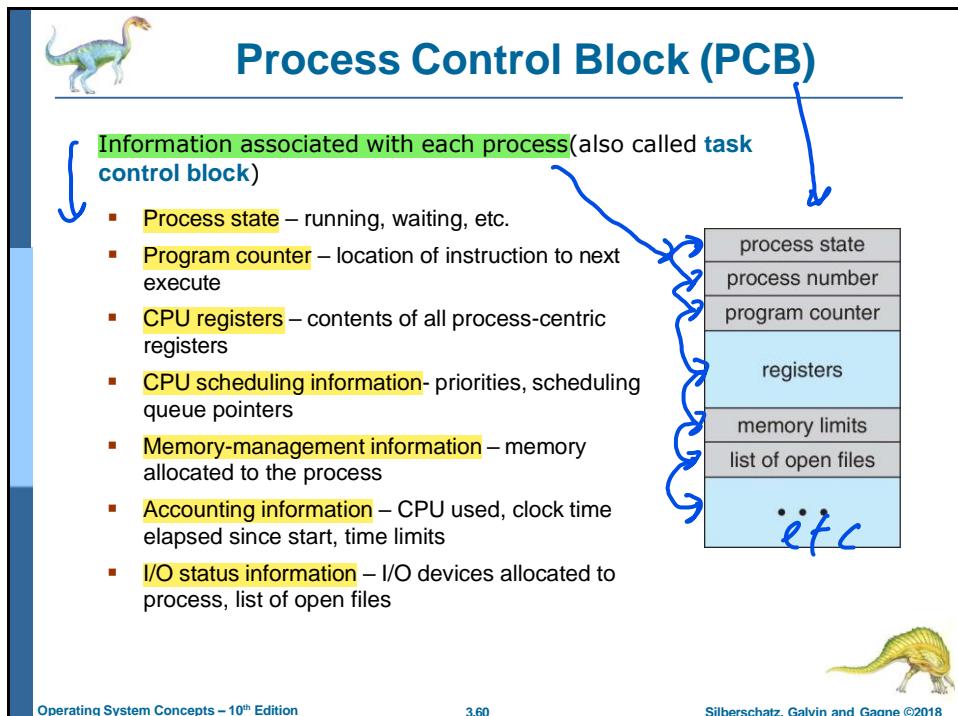
Silberschatz, Galvin and Gagne ©2018

58

29



59



60

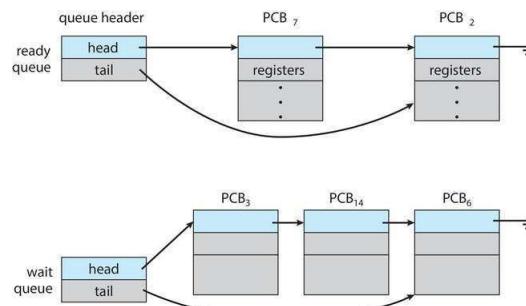


## Threads (Self Study)

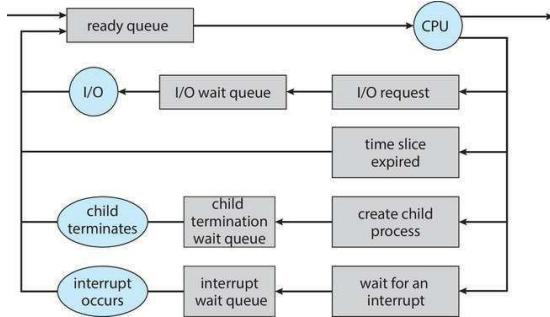
- So far, process has a single thread of execution
- Consider having multiple program counters per process
  - Multiple locations can execute at once
  - Multiple threads of control -> **threads**
- Must then have storage for thread details, multiple program counters in PCB
- Explore in detail in Chapter 4



## Ready and Wait Queues (Self Study)



## Representation of Process Scheduling (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

3.63

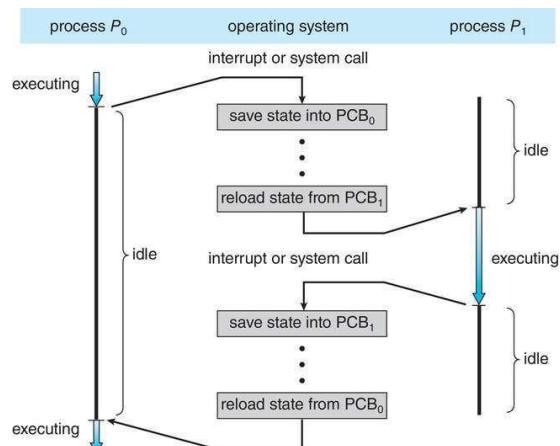
Silberschatz, Galvin and Gagne ©2018

63



## CPU Switch From Process to Process

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



Operating System Concepts – 10<sup>th</sup> Edition

3.64

Silberschatz, Galvin and Gagne ©2018

64



## Context Switch

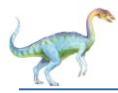
- When CPU switches to another process, the system must **save the state** of the old process and load the **saved state** for the new process via a **context switch**
- **Context** of a process represented in the PCB
- Context-switch time is pure overhead; the system does no useful work while switching
  - 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





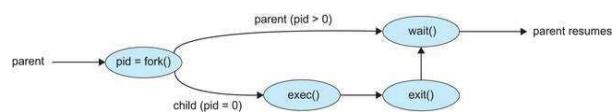
## Process Creation

- 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
  - Parent and children share all resources
  - Children share subset of parent's resources
  - Parent and child share no resources
- Execution options
  - Parent and children execute concurrently
  - Parent waits until children terminate



## Process Creation (Cont.)

- Address space
  - Child duplicate of parent
  - Child has a program loaded into it
- UNIX examples
  - `fork()` system call creates new process
  - `exec()` system call used after a `fork()` to replace the process' memory space with a new program
  - Parent process calls `wait()` waiting for the child to terminate





## 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



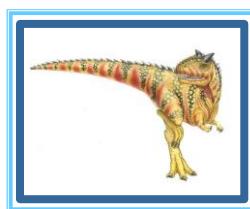
## Process Termination

- 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 – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

71



## Outline

- Overview
- Multicore Programming
- Multithreading Models
- Thread Libraries
- Threading Issues

Operating System Concepts – 10<sup>th</sup> Edition

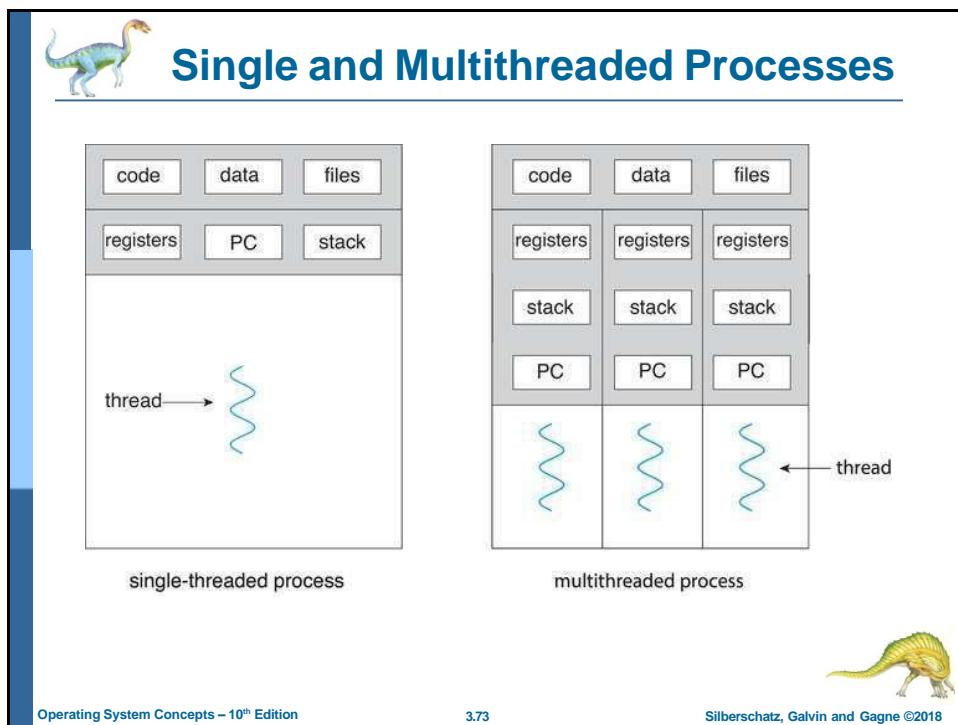
3.72

Silberschatz, Galvin and Gagne ©2018

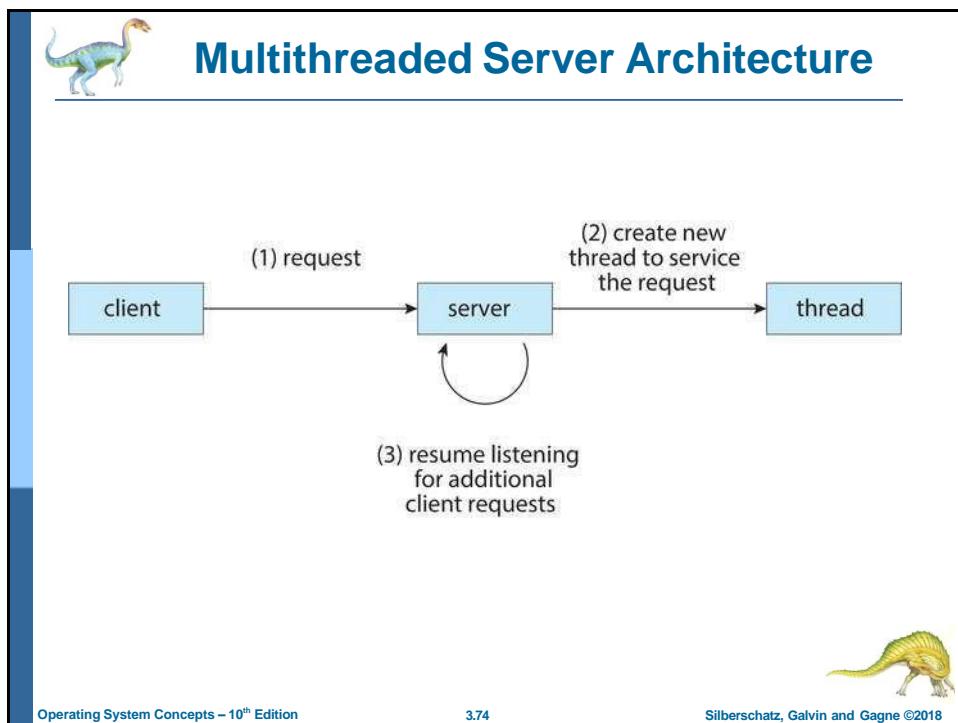


72

36



73



74



## 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



## Multicore Programming

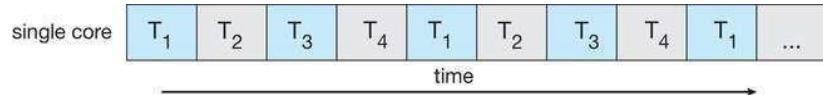
- **Multicore** or **multiprocessor** systems puts pressure on programmers, challenges include:
  - **Dividing activities**
  - **Balance**
  - **Data splitting**
  - **Data dependency**
  - **Testing and debugging**
- **Parallelism** implies a system can perform more than one task simultaneously
- **Concurrency**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



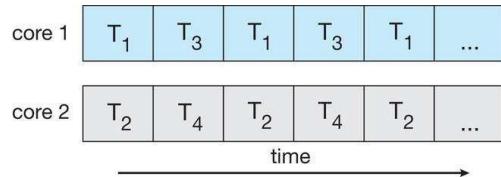


## Concurrency vs. Parallelism

- Concurrent execution on single-core system:



- Parallelism on a multi-core system:



## Multicore Programming (Self Study)

- Types of parallelism
  - Data parallelism** – distributes subsets of the same data across multiple cores, same operation on each
  - Task parallelism** – distributing threads across cores, each thread performing unique operation



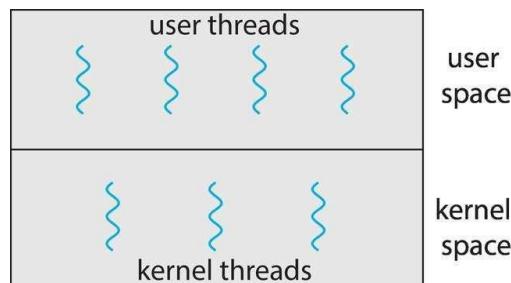


## User Threads and Kernel Threads

- **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





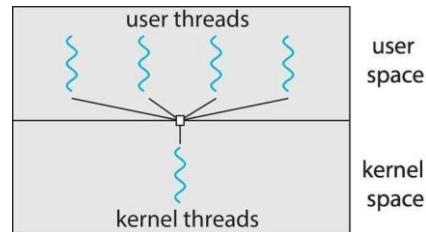
## Multithreading Models

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



## Many-to-One

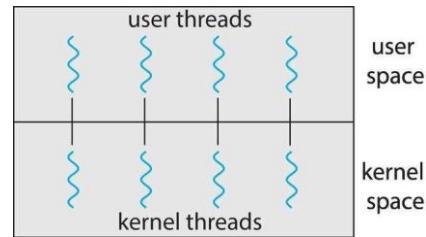
- Many user-level threads mapped to single kernel thread
- One thread blocking causes all to block
- Multiple threads may not run in parallel on multicore system because only one may be in kernel at a time
- Few systems currently use this model
- Examples:
  - Solaris Green Threads
  - GNU Portable Threads





## 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 – 10<sup>th</sup> Edition

3.83

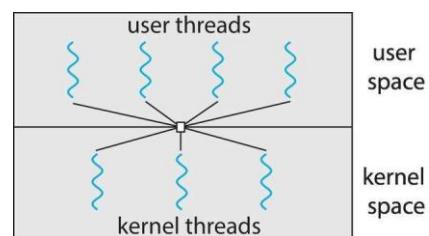
Silberschatz, Galvin and Gagne ©2018

83



## 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 – 10<sup>th</sup> Edition

3.84

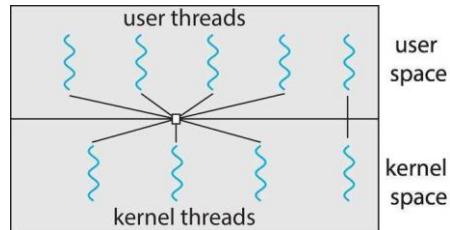
Silberschatz, Galvin and Gagne ©2018

84



## 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





## Pthreads (Self Study)

- May be provided either as user-level or kernel-level
- A POSIX standard (IEEE 1003.1c) API for thread creation and synchronization
- **Specification**, not *implementation*
- API specifies behavior of the thread library, implementation is up to development of the library
- Common in UNIX operating systems (Linux & Mac OS X)



## Pthreads Example (Self Study)

```
#include <pthread.h>
#include <stdio.h>

#include <stdlib.h>

int sum; /* this data is shared by the thread(s) */
void *runner(void *param); /* threads call this function */

int main(int argc, char *argv[])
{
    pthread_t tid; /* the thread identifier */
    pthread_attr_t attr; /* set of thread attributes */

    /* set the default attributes of the thread */
    pthread_attr_init(&attr);
    /* create the thread */
    pthread_create(&tid, &attr, runner, argv[1]);
    /* wait for the thread to exit */
    pthread_join(tid,NULL);

    printf("sum = %d\n",sum);
}
```





## Java Threads (Self Study)

- Java threads are managed by the JVM
- Typically implemented using the threads model provided by underlying OS
- Java threads may be created by:
  - Extending Thread class
  - Implementing the Runnable interface

```
public interface Runnable  
{  
    public abstract void run();  
}
```

- Standard practice is to implement Runnable interface



## Java Threads (Self Study)

### Implementing Runnable interface:

```
class Task implements Runnable  
{  
    public void run() {  
        System.out.println("I am a thread.");  
    }  
}
```

### Creating a thread:

```
Thread worker = new Thread(new Task());  
worker.start();
```

### Waiting on a thread:

```
try {  
    worker.join();  
}  
catch (InterruptedException ie) { }
```





## Threading Issues (Self Study)

- Semantics of **fork()** and **exec()** system calls
- Signal handling
  - Synchronous and asynchronous
- Thread cancellation of target thread
  - Asynchronous or deferred
- Thread-local storage
- Scheduler Activations



## Semantics of fork() and exec() (Self Study)

- Does **fork()** duplicate only the calling thread or all threads?
  - Some UNIXes have two versions of fork
- **exec()** usually works as normal – replace the running process including all threads





## 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 canceled
- 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
- | 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





## Thread Cancellation in Java (Self Study)

- Deferred cancellation uses the `interrupt()` method, which sets the interrupted status of a thread.

```
Thread worker;
```

```
    . . .
```

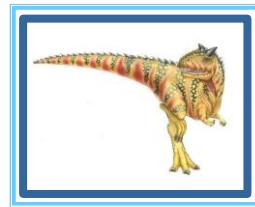
```
/* set the interruption status of the thread */
worker.interrupt()
```

- A thread can then check to see if it has been interrupted:

```
while (!Thread.currentThread().isInterrupted()) {
    . . .
}
```

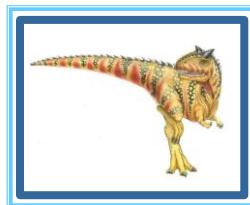


## End of Chapter 2



12 / 11

# Chapter 3: CPU Scheduling



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

97



## Outline

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

Operating System Concepts – 10<sup>th</sup> Edition

5.98

Silberschatz, Galvin and Gagne ©2018

98

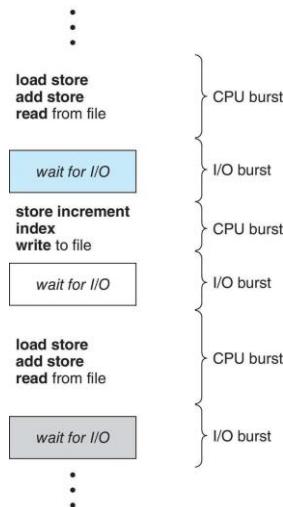


49



## 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 – 10<sup>th</sup> Edition

5.99

Silberschatz, Galvin and Gagne ©2018

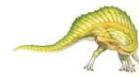
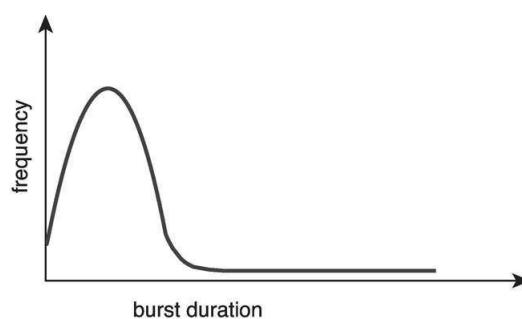
99



## Histogram of CPU-burst Times

Large number of short bursts

Small number of longer bursts



Operating System Concepts – 10<sup>th</sup> Edition

5.100

Silberschatz, Galvin and Gagne ©2018

100



## CPU Scheduler

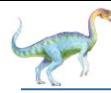
- The **CPU scheduler** selects from among the processes in ready queue, and allocates a CPU core to one of them
  - Queue may be ordered in various ways
- CPU scheduling decisions may take place when a process:
  1. Switches from running to waiting state
  2. Switches from running to ready state
  3. 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.





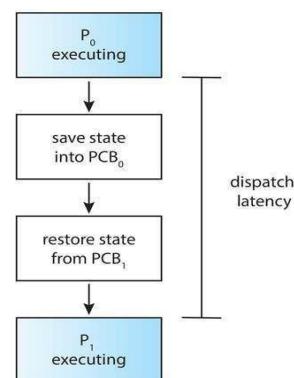
## Preemptive Scheduling and Race Conditions

- 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.



## 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





## Scheduling Criteria

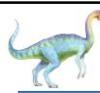
- **CPU utilization** – keep the CPU as busy as possible
- **Throughput** – # of processes that complete their execution per time unit
- **Turnaround time** – amount of time to execute a particular process
- **Waiting time** – amount of time a process has been waiting in the ready queue
- **Response 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

Silberschatz, Galvin and Gagne ©2018

Operating System Concepts – 10<sup>th</sup> Edition

5.107



107



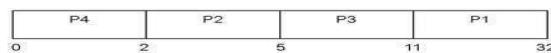
## Shortest-Job-First (SJF) Scheduling

### 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 =  $(11+2+5+0) / 4 = 4.5 \text{ ms}$



Operating System Concepts – 10<sup>th</sup> Edition

5.108

Silberschatz, Galvin and Gagne ©2018

108



## 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 | P1 |
| 0  | 1  | 4  | 6  | 12 | 32 |

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

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 – 10<sup>th</sup> Edition      5.109

109



## 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                          | $= (15-2) \rightarrow 13$                   |
| <b>Total Waiting Time</b>   | <b>=44</b>                                  |
| <b>Average Waiting Time</b> | <b>=44/4 <math>\rightarrow</math> 11 ms</b> |



Operating System Concepts – 10<sup>th</sup> Edition      5.110  
Silberschatz, Galvin and Gagne ©2018

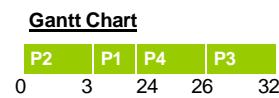
110



## 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 =  $53/4 = 13.25$ ms

| Process | Waiting Time |
|---------|--------------|
| P1      | 3            |
| P2      | 0            |
| P3      | 26           |
| P4      | 24           |

Operating System Concepts – 10<sup>th</sup> Edition

5.111

Silberschatz, Galvin and Gagne ©2018



111

12 / 13



## 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 – 10<sup>th</sup> Edition

5.112

Silberschatz, Galvin and Gagne ©2018

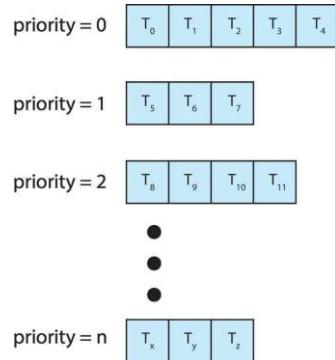


112



## Multilevel Queue

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



Operating System Concepts – 10<sup>th</sup> Edition

5.113

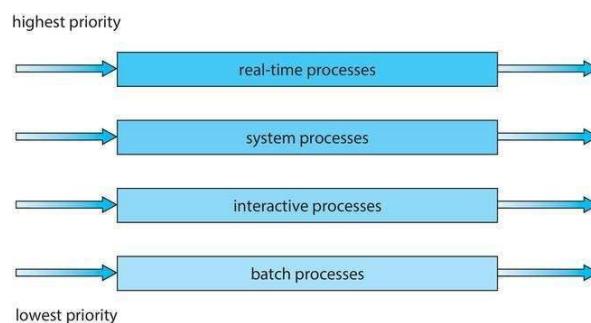
Silberschatz, Galvin and Gagne ©2018

113



## Multilevel Queue

- Prioritization based upon process type



Operating System Concepts – 10<sup>th</sup> Edition

5.114

Silberschatz, Galvin and Gagne ©2018

114



## Multilevel Feedback Queue (Self Study)

- A process can move between the various queues.
- Multilevel-feedback-queue scheduler defined by the following parameters:
  - Number of queues
  - Scheduling algorithms for each queue
  - Method used to determine when to upgrade a process
  - Method used to determine when to demote a process
  - Method used to determine which queue a process will enter when that process needs service
- Aging can be implemented using multilevel feedback queue



Operating System Concepts – 10<sup>th</sup> Edition

5.115

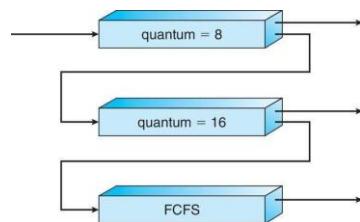
Silberschatz, Galvin and Gagne ©2018

115



## (Self Study)

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

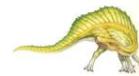


Operating System Concepts – 10<sup>th</sup> Edition

5.116

Silberschatz, Galvin and Gagne ©2018

116





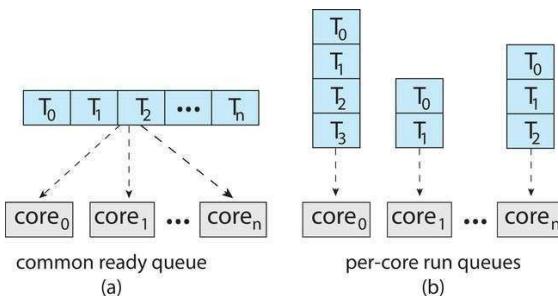
## Multiple-Processor Scheduling

- 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

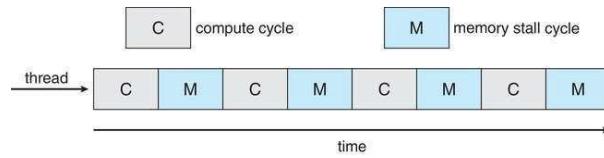
- Symmetric multiprocessing (SMP) is where each processor is self scheduling.
- All threads may be in a common ready queue (a)
- Each processor may have its own private queue of threads (b)





## 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 – 10<sup>th</sup> Edition

5.119

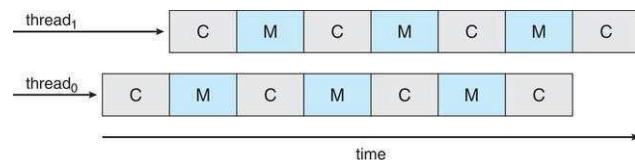
Silberschatz, Galvin and Gagne ©2018

119



## Study)

- Each core has > 1 hardware threads.
- If one thread has a memory stall, switch to another thread!
- Figure



Operating System Concepts – 10<sup>th</sup> Edition

5.120

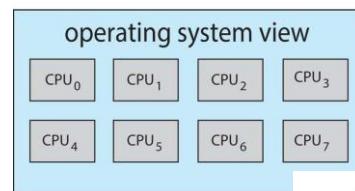
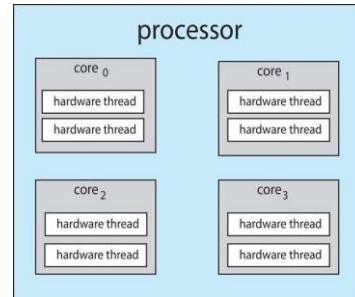
Silberschatz, Galvin and Gagne ©2018

120



## Study)

- **Chip-multithreading (CMT)** assigns each core multiple hardware threads. (Intel refers to this as **hyperthreading**.)
- On a quad-core system with 2 hardware threads per core, the operating system sees 8 logical processors.



Operating System Concepts – 10<sup>th</sup> Edition

5.121

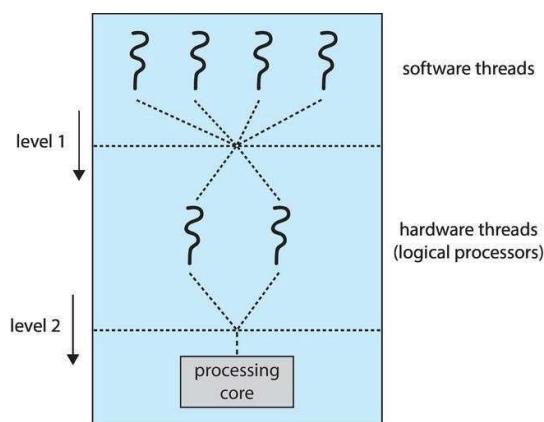
Silberschatz, Galvin and Gagne ©2018

121



## Study)

- Two levels of scheduling:
  1. The operating system deciding which software thread to run on a logical CPU
  2. How each core decides which hardware thread to run on the physical core.



Operating System Concepts – 10<sup>th</sup> Edition

5.122

Silberschatz, Galvin and Gagne ©2018

122



## Multiple-Processor Scheduling – Load Balancing (Self Study)



- If SMP, need to keep all CPUs loaded for efficiency
- **Load balancing** attempts to keep workload evenly distributed
- **Push migration** – periodic task checks load on each processor, and if found pushes task from overloaded CPU to other CPUs
- **Pull migration** – idle processors pulls waiting task from busy processor



## Multiple-Processor Scheduling – Processor Affinity (Self Study)



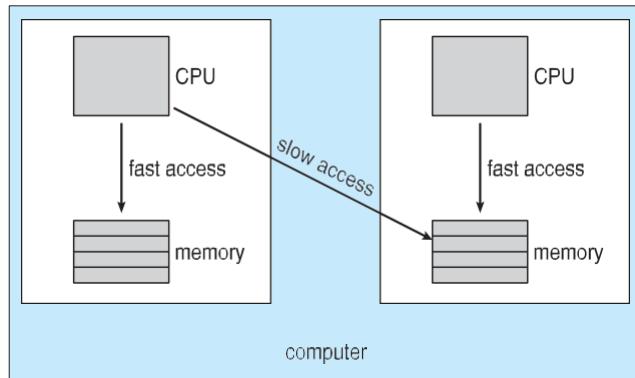
- When a thread has been running on one processor, the cache contents of that processor stores the memory accesses by that thread.
- We refer to this as a thread having affinity for a processor (i.e., “processor affinity”)
- Load balancing may affect processor affinity as a thread may be moved from one processor to another to balance loads, yet that thread loses the contents of what it had in the cache of the processor it was moved off of.
- **Soft affinity** – the operating system attempts to keep a thread running on the same processor, but no guarantees.
- **Hard affinity** – allows a process to specify a set of processors it may run on.





## (Self Study)

If the operating system is **NUMA-aware**, it will assign memory closer to the CPU the thread is running on.



## Real-Time CPU Scheduling

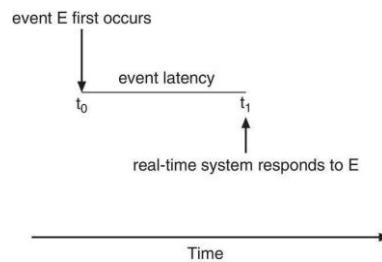
- Can present obvious challenges
- **Soft real-time systems** – Critical real-time tasks have the highest priority, but no guarantee as to when tasks will be scheduled
- **Hard real-time systems** – task must be serviced by its deadline





## 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 – 10<sup>th</sup> Edition

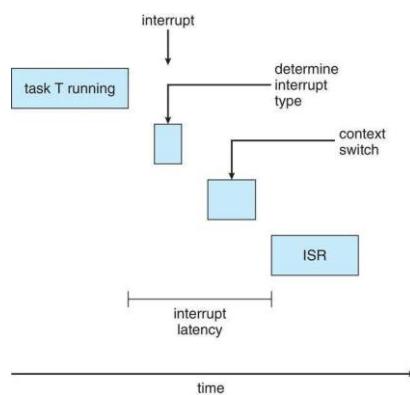
5.127

Silberschatz, Galvin and Gagne ©2018

127



## Interrupt Latency (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

5.128

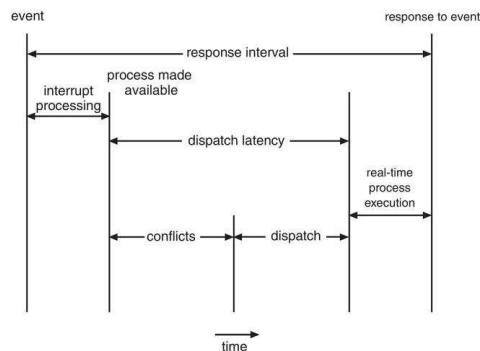
Silberschatz, Galvin and Gagne ©2018

128



## Dispatch Latency (Self Study)

- Conflict phase of dispatch latency:
  1. Preemption of any process running in kernel mode
  2. Release by low-priority process of resources needed by high-priority processes



Operating System Concepts – 10<sup>th</sup> Edition

5.129

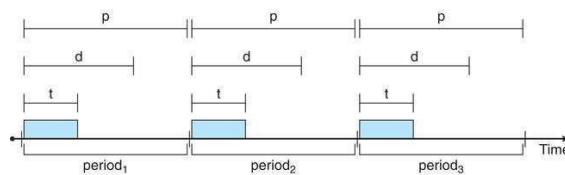
Silberschatz, Galvin and Gagne ©2018

129



## Priority-based Scheduling (Self Study)

- For real-time scheduling, scheduler must support preemptive, priority-based scheduling
  - But only guarantees soft real-time
- For hard real-time must also provide ability to meet deadlines
- Processes have new characteristics: **periodic** ones require CPU at constant intervals
  - Has processing time  $t$ , deadline  $d$ , period  $p$
  - $0 \leq t \leq d \leq p$
  - **Rate** of periodic task is  $1/p$



Operating System Concepts – 10<sup>th</sup> Edition

5.130

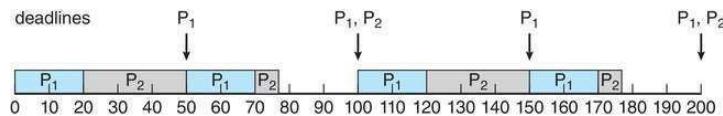
Silberschatz, Galvin and Gagne ©2018

130



## Study)

- A priority is assigned based on the inverse of its period
- Shorter periods = higher priority;
- Longer periods = lower priority
- $P_1$  is assigned a higher priority than  $P_2$ .



Operating System Concepts – 10<sup>th</sup> Edition

5.131

Silberschatz, Galvin and Gagne ©2018

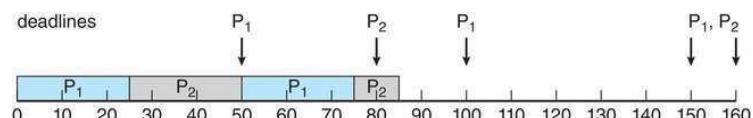


131



## Missed Deadlines with Rate Monotonic Scheduling (Self Study)

- Process  $P_2$  misses finishing its deadline at time 80
- Figure



Operating System Concepts – 10<sup>th</sup> Edition

5.132

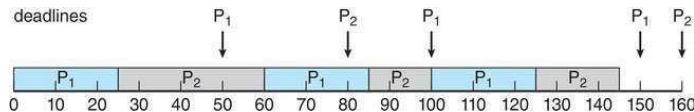
Silberschatz, Galvin and Gagne ©2018



132

## Earliest Deadline First Scheduling (EDF) (Self Study)

- Priorities are assigned according to deadlines:
  - The earlier the deadline, the higher the priority
  - The later the deadline, the lower the priority
- Figure



Operating System Concepts – 10<sup>th</sup> Edition

5.133

Silberschatz, Galvin and Gagne ©2018



133

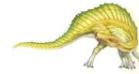
## (Self Study)

- $T$  shares are allocated among all processes in the system
- An application receives  $N$  shares where  $N < T$
- This ensures each application will receive  $N / T$  of the total processor time

Operating System Concepts – 10<sup>th</sup> Edition

5.134

Silberschatz, Galvin and Gagne ©2018



134

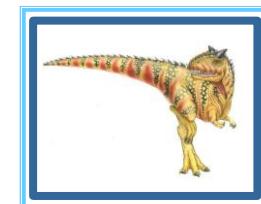


## (Self Study)

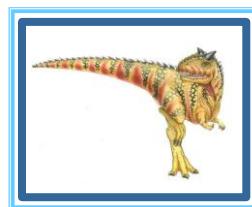
- The POSIX.1b standard
- API provides functions for managing real-time threads
- Defines two scheduling classes for real-time threads:
  1. `SCHED_FIFO` - threads are scheduled using a FCFS strategy with a FIFO queue. There is no time-slicing for threads of equal priority
  2. `SCHED_RR` - similar to `SCHED_FIFO` except time-slicing occurs for threads of equal priority
- Defines two functions for getting and setting scheduling policy:
  1. `pthread_attr_getsched_policy(pthread_attr_t *attr, int *policy)`
  2. `pthread_attr_setsched_policy(pthread_attr_t *attr, int policy)`



## End of Chapter 3



# Chapter 4: Synchronization & Deadlocks



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

137



## Outline

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

Operating System Concepts – 10<sup>th</sup> Edition

6.138

Silberschatz, Galvin and Gagne ©2018

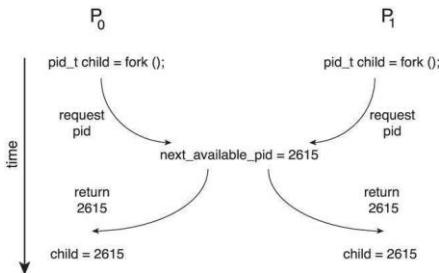
138





## 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!



Operating System Concepts – 10<sup>th</sup> Edition

6.139

Silberschatz, Galvin and Gagne ©2018

139



## 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**



Operating System Concepts – 10<sup>th</sup> Edition

6.140

Silberschatz, Galvin and Gagne ©2018

140



## 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





## Semaphore Usage Example (Self Study)

- Solution to the CS Problem
  - Create a semaphore “`mutex`” initialized to 1

```
    wait(mutex);
```

CS

```
    signal(mutex);
```

- Consider  $P_1$  and  $P_2$  that with two statements  $S_1$ , and  $S_2$  and the requirement that  $S_1$  to happen before  $S_2$

- Create a semaphore “`synch`” initialized to 0

$P_1$ :

```
    S1;
```

```
    signal(synch);
```

$P_2$ :

```
    wait(synch);
```

```
    S2;
```



## Semaphore Implementation (Self Study)

- Must guarantee that no two processes can execute the `wait()` and `signal()` on the same semaphore at the same time
- Thus, the implementation becomes the critical section problem where the `wait` and `signal` code are placed in the critical section
- Could now have **busy waiting** in critical section implementation
  - But implementation code is short
  - Little busy waiting if critical section rarely occupied
- Note that applications may spend lots of time in critical sections and therefore this is not a good solution





## Study)

- Incorrect use of semaphore operations:
  - `signal(mutex) ... wait(mutex)`
  - `wait(mutex) ... wait(mutex)`
  - Omitting of `wait (mutex)` and/or `signal (mutex)`
- These – and others – are examples of what can occur when semaphores and other synchronization tools are used incorrectly.



## Monitors (Self Study)

- A high-level abstraction that provides a convenient and effective mechanism for process synchronization
- *Abstract data type*, internal variables only accessible by code within the procedure
- Only one process may be active within the monitor at a time
- Pseudocode syntax of a monitor:

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

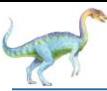
    procedure P2 (...) { .... }

    procedure Pn (...) {.....}

    initialization code (...) { ... }
}
```



## Monitor Implementation Using Semaphores (Self Study)



- Variables

```
semaphore mutex
mutex = 1
```

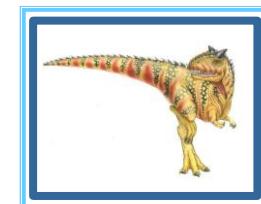
- Each procedure *P* is replaced by

```
wait(mutex);
...
body of P;
...
signal(mutex);
```

- Mutual exclusion within a monitor is ensured



## Chapter 4 (Cont'd): Synchronization Examples





## Outline

- Explain the bounded-buffer synchronization problem
- Explain the readers-writers synchronization problem
- Explain and dining-philosophers synchronization problems



## 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



## Bounded Buffer Problem (Cont.) (Self Study)

- The structure of the producer process

```
while (true) {  
    ...  
    /* produce an item in next_produced */  
    ...  
    wait(empty);  
    wait(mutex);  
    ...  
    /* add next_produced to the buffer */  
    ...  
    signal(mutex);  
    signal(full);  
}
```



## Bounded Buffer Problem (Cont.)

### (Self Study)

- The structure of the consumer process

```
while (true) {  
    wait(full);  
    wait(mutex);  
  
    ...  
    /* remove an item from buffer to next_consumed */  
    ...  
    signal(mutex);  
    signal(empty);  
  
    ...  
    /* consume the item in next consumed */  
    ...  
}
```

Operating System Concepts – 10<sup>th</sup> Edition

6.155

Silberschatz, Galvin and Gagne ©2018



155

## 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 – 10<sup>th</sup> Edition

6.156

Silberschatz, Galvin and Gagne ©2018



156



## 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



## Readers-writers Problem (Cont.) (Self Study)

- The structure of a writer process

```
while (true) {
    wait(rw_mutex);
    ...
    /* writing is performed */
    ...
    signal(rw_mutex);
}
```



## Readers-Writers Problem (Cont.) (Self Study)

- The structure of a reader process

```
while (true){  
    wait(mutex);  
    read_count++;  
    if (read_count == 1) /* first reader */  
        wait(rw_mutex);  
    signal(mutex);  
    ...  
    /* reading is performed */  
    ...  
    wait(mutex);  
    read_count--;  
    if (read_count == 0) /* last reader */  
        signal(rw_mutex);  
    signal(mutex);  
}
```



Operating System Concepts – 10<sup>th</sup> Edition

6.159

Silberschatz, Galvin and Gagne ©2018

159



## Dining-Philosophers Problem

- N philosophers sit at a round table with a bowl 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



Operating System Concepts – 10<sup>th</sup> Edition

6.160

Silberschatz, Galvin and Gagne ©2018

160



## Dining-Philosophers Problem Algorithm (Self Study)

- Semaphore Solution

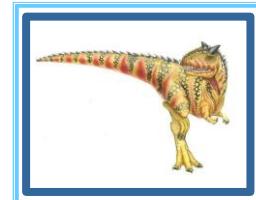
- The structure of Philosopher *i*:

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

- What is the problem with this algorithm?



## Chapter 4 (Cont'd): Deadlocks





## Outline

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



## 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





## Study)

- Data:
  - A semaphore  $s_1$  initialized to 1
  - A semaphore  $s_2$  initialized to 1
- Two threads  $T_1$  and  $T_2$
- $T_1$ :

```
wait(s1)  
wait(s2)
```
- $T_2$ :

```
wait(s2)  
wait(s1)
```



## 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

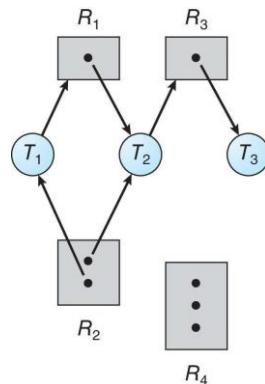
A set of vertices  $V$  and a set of edges  $E$ .

- $V$  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  $R_j \rightarrow T_i$



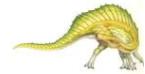
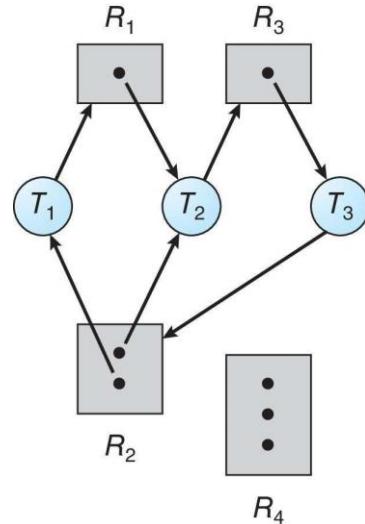
## Resource Allocation Graph Example

- One instance of  $R_1$
- Two instances of  $R_2$
- One instance of  $R_3$
- Three instances 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  $R_3$
- $T_3$  holds one instance of  $R_3$





## Resource Allocation Graph with a Deadlock



Operating System Concepts – 10<sup>th</sup> Edition

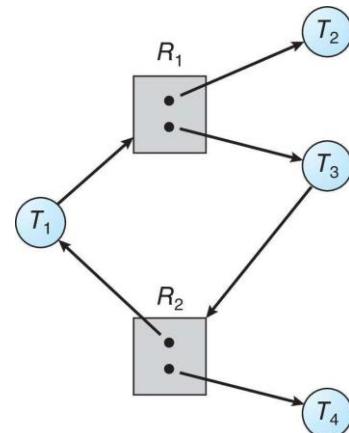
6.169

Silberschatz, Galvin and Gagne ©2018

169



## Graph with a Cycle But no Deadlock



Operating System Concepts – 10<sup>th</sup> Edition

6.170

Silberschatz, Galvin and Gagne ©2018

170



## 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  $\langle T_1, T_2, \dots, T_r \rangle$  of all the threads in the system 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



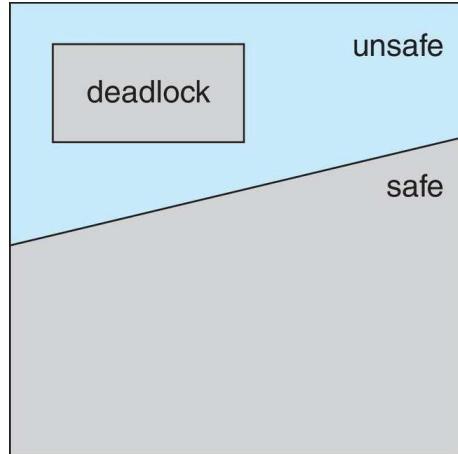
## Basic Facts

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





## Safe, Unsafe, Deadlock State



Operating System Concepts – 10<sup>th</sup> Edition

6.179

Silberschatz, Galvin and Gagne ©2018

179



## 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 – 10<sup>th</sup> Edition

6.180

Silberschatz, Galvin and Gagne ©2018

180

90



## 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 – 10<sup>th</sup> Edition

6.181

Silberschatz, Galvin and Gagne ©2018

181



### 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

Silberschatz, Galvin and Gagne ©2018

Operating System Concepts – 10<sup>th</sup> Edition

6.182

182

## Data Structures for the Banker's Algorithm (Self Study)



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

- **Available:** Vector of length  $m$ . If  $\text{Available}[j] = k$ , there are  $k$  instances of resource type  $R_j$  available
- **Max:**  $n \times m$  matrix. If  $\text{Max}[i,j] = k$ , then process  $T_i$  may request at most  $k$  instances of resource type  $R_j$
- **Allocation:**  $n \times m$  matrix. If  $\text{Allocation}[i,j] = k$  then  $T_i$  is currently allocated  $k$  instances of  $R_j$
- **Need:**  $n \times m$  matrix. If  $\text{Need}[i,j] = k$ , then  $T_i$  may need  $k$  more instances of  $R_j$  to complete its task

$$\text{Need}[i,j] = \text{Max}[i,j] - \text{Allocation}[i,j]$$



## Safety Algorithm (Self Study)



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

**Work = Available**

**Finish[1] = raise** for  $i = 0, 1, \dots, n-1$

2. Find an  $i$  such that both:
  - (a) **Finish[i] = false**
  - (b) **Need<sub>i</sub> ⊑ Work**  
If no such  $i$  exists, go to step 4
3. **Work = Work + Allocation<sub>i</sub>**,  
**Finish[i] = true**  
go to step 2
4. If **Finish[i] == true** for all  $i$ , then the system is in a safe state



## Resource-Request Algorithm for Process $P_i$ (Self Study)

$\text{Request}_i$  = request vector for process  $T_i$ . If  $\text{Request}_i[j] = k$  then process  $T_i$  wants  $k$  instances of resource type  $R_j$

1. If  $\text{Request}_i \geq \text{Need}_i$ , go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim
2. If  $\text{Request}_i \leq \text{Available}$ , go to step 3. Otherwise  $T_i$  must wait, since resources are not available
3. Pretend to allocate requested resources to  $T_i$  by modifying the state as follows:

$$\text{Available} = \text{Available} - \text{Request}_i;$$

$$\text{Allocation}_i = \text{Allocation}_i + \text{Request}_i;$$

$$\text{Need}_i = \text{Need}_i - \text{Request}_i;$$

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



## Deadlock Detection

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



## Single Instance of Each Resource Type (Self Study)

- Maintain **wait-for** graph
  - Nodes are threads
  - $T_i \sqsupseteq T_j$  if  $T_i$  is waiting for  $T_j$
- Periodically invoke an algorithm that searches for a cycle in the graph. If there is a cycle, there exists a deadlock
- An algorithm to detect a cycle in a graph requires an order of  $n^2$  operations, where  $n$  is the number of vertices in the graph

Operating System Concepts – 10<sup>th</sup> Edition

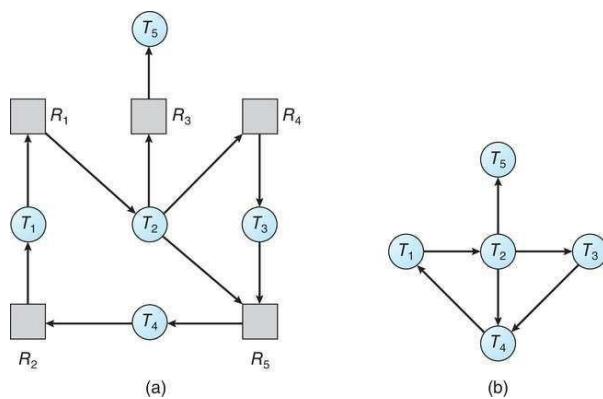
6.187

Silberschatz, Galvin and Gagne ©2018



187

## Resource-Allocation Graph and Wait-for Graph (Self Study)



Resource-Allocation Graph

Corresponding wait-for graph

Operating System Concepts – 10<sup>th</sup> Edition

6.188

Silberschatz, Galvin and Gagne ©2018



188



## (Self Study)

- **Available:** A vector of length  $m$  indicates the number of available resources of each type
- **Allocation:** An  $n \times m$  matrix defines the number of resources of each type currently allocated to each thread.
- **Request:** An  $n \times m$  matrix indicates the current request of each thread. If  $\text{Request}[i][j] = k$ , then thread  $T_i$  is requesting  $k$  more instances of resource type  $R_j$ .



## Detection Algorithm (Self Study)

1. Let  $Work$  and  $Finish$  be vectors of length  $m$  and  $n$ , respectively  
Initialize:
  - a)  $Work = Available$
  - b) For  $i = 1, 2, \dots, n$ , if  $Allocation_{i, \cdot} \neq 0$ , then  $Finish[i] = \text{false}$ ; otherwise,  $Finish[i] = \text{true}$
2. Find an index  $i$  such that both:
  - a)  $Finish[i] == \text{false}$
  - b)  $\text{Request}_{i, \cdot} \leq Work$

If no such  $i$  exists, go to step 4



## Detection Algorithm (Cont.) (Self Study)



3.  $Work = Work + Allocation$ ,  
 $Finish[i] = true$   
go to step 2
4. If  $Finish[i] = false$ , for some  $i$ ,  $1 \leq i \leq n$ , then the system is in deadlock state. Moreover, if  $Finish[i] == false$ , then  $T_i$  is deadlocked

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



## 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  $Finish[i] = true$  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?

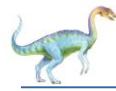
- Can reclaim resources held by thread  $T_0$ , but insufficient resources to fulfill other processes' requests
- Deadlock exists, consisting of processes  $T_1$ ,  $T_2$ ,  $T_3$ , and  $T_4$



## Detection-Algorithm Usage (Self Study)

- When, and how often, to invoke depends on:
  - How often a deadlock is likely to occur?
  - How many processes will need to be rolled back?
    - ▶ one for each disjoint cycle
- If detection algorithm is invoked arbitrarily, there may be many cycles in the resource graph and so we would not be able to tell which of the many deadlocked threads "caused" the deadlock.





## 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?

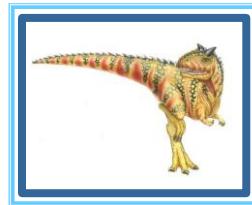


## Recovery from Deadlock: Resource Preemption

- **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



## End of Chapter 4

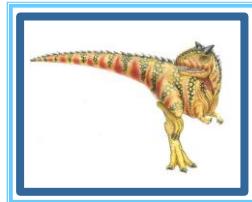


Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

197

## Chapter 5: Memory Management (Main Memory)



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

198

99



## Chapter 5: Memory Management

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



Operating System Concepts – 10<sup>th</sup> Edition

9.199

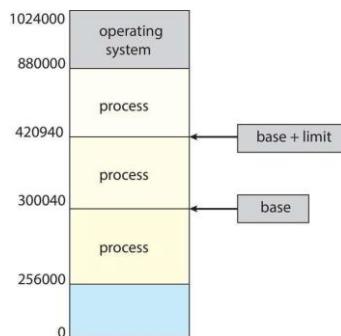
Silberschatz, Galvin and Gagne ©2018

199



## 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 – 10<sup>th</sup> Edition

9.200

Silberschatz, Galvin and Gagne ©2018

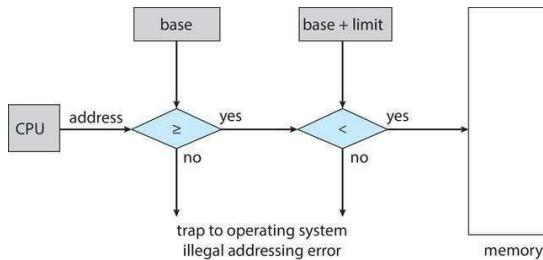
200

100



## 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

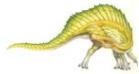
- Programs on disk, ready to be brought into memory to execute from an **input queue**
  - Without support, must be loaded into address 0000
- Inconvenient to have first user process physical address always at 0000
  - How can it not be?
- Addresses represented in different ways at different stages of 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)



## Logical vs. Physical Address Space

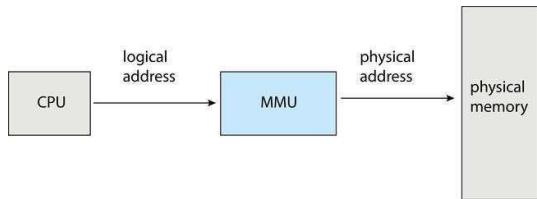
- The concept of a logical address space that is bound to a separate **physical address space** is central to proper memory management
  - **Logical address** – generated by the CPU; also referred to as **virtual address**
  - **Physical address** – address seen by the memory unit
- Logical and physical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





## Memory-Management Unit (MMU)

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



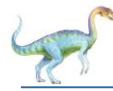
- Many methods possible, covered in the rest of this chapter



## Memory-Management Unit (Cont.) (Self Study)

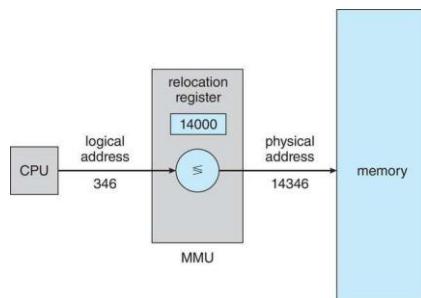
- Consider simple scheme, which is a generalization of the base-register scheme.
- The base register now called **relocation register**
- The value in the relocation register is added to every address generated by a user process at the time it is sent to memory
- The user program deals with *logical* addresses; it never sees the *real* physical addresses
  - Execution-time binding occurs when reference is made to location in memory
  - Logical address bound to physical addresses





## 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 – 10<sup>th</sup> Edition

9.207

Silberschatz, Galvin and Gagne ©2018

207



## 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 – 10<sup>th</sup> Edition

9.208

Silberschatz, Galvin and Gagne ©2018

208

104



## 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



## Contiguous Allocation (Self Study)

- Main memory must support both OS and user processes
- Limited resource, must allocate efficiently
- Contiguous allocation is one early method
- Main memory usually into two **partitions**:
  - Resident operating system, usually held in low memory with interrupt vector
  - User processes then held in high memory
  - Each process contained in single contiguous section of memory



## Contiguous Allocation (Cont.) (Self Study)

- Relocation registers used to protect user processes from each other, and from changing operating-system code and data
  - Base register contains value of smallest physical address
  - Limit register contains range or logical addresses – each logical address must be less than the limit register
  - MMU maps logical address *dynamically*
  - Can then allow actions such as kernel code being **transient** and kernel changing size

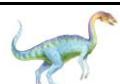
Operating System Concepts – 10<sup>th</sup> Edition

9.211

Silberschatz, Galvin and Gagne ©2018

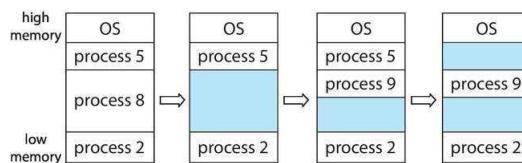


211



## 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 – 10<sup>th</sup> Edition

9.212

Silberschatz, Galvin and Gagne ©2018



212



## Dynamic Storage-Allocation Problem

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

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

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



## Fragmentation

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





## Fragmentation (Cont.) (Self Study)

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



## Paging

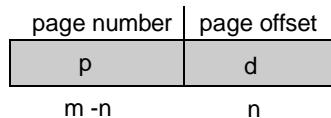
- Physical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available
  - Avoids external fragmentation
  - Avoids problem of varying sized memory chunks
- Divide physical memory into fixed-sized blocks called **frames**
  - Size is power of 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

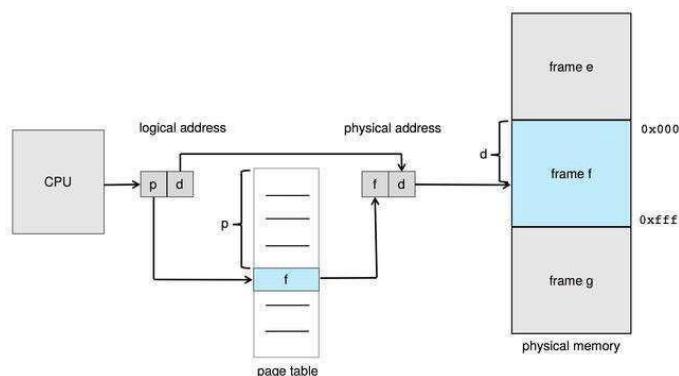
- Address generated by CPU is divided into:
  - **Page number (*p*)** – used as an index into a **page table** which contains base address of each page in physical memory
  - **Page offset (*d*)** – combined with base address to define the physical memory address that is sent to the memory unit

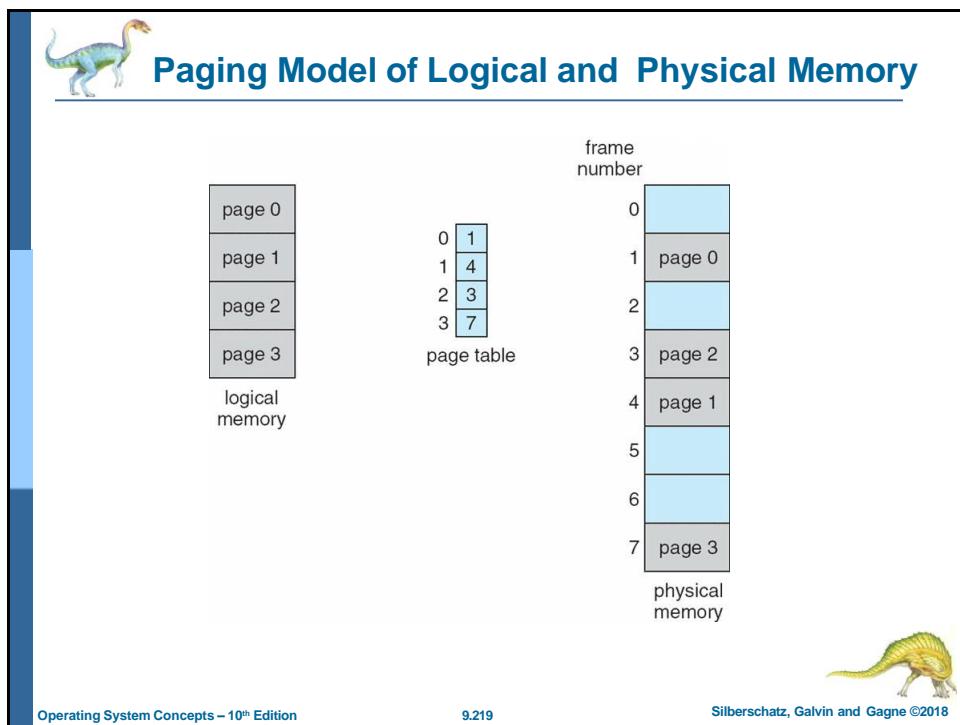


- For given logical address space  $2^m$  and page size  $2^n$

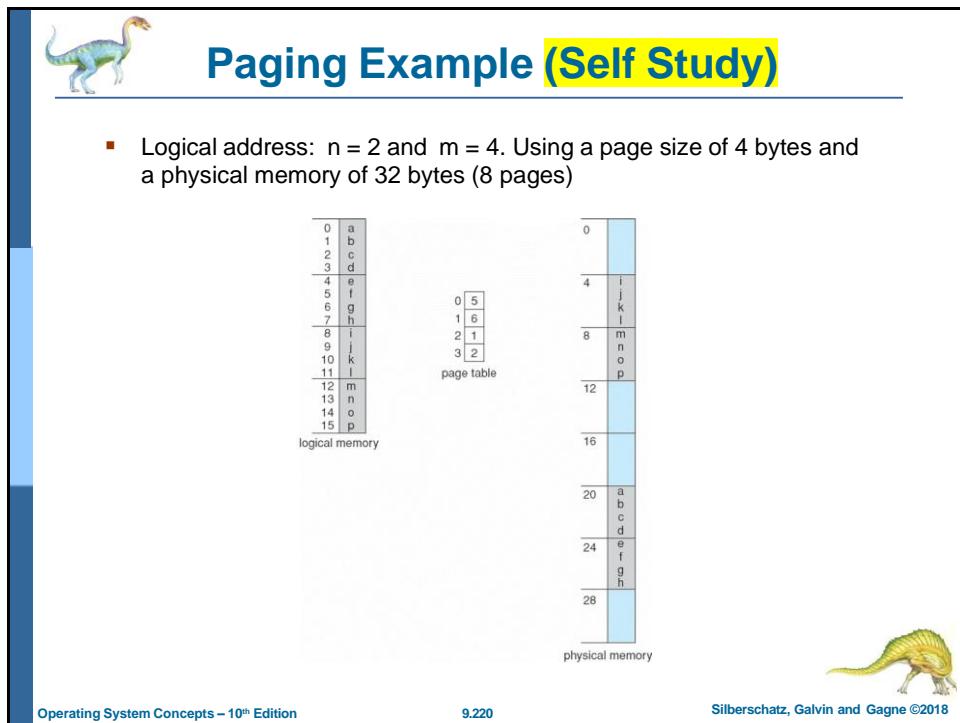


## Paging Hardware (Self Study)





219



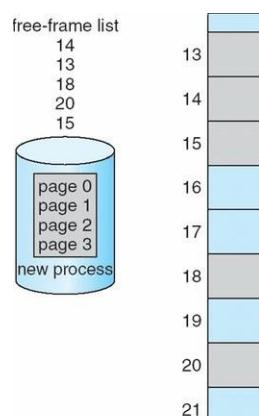
220

## Paging -- Calculating internal fragmentation (Self Study)

- Page size = 2,048 bytes
- Process size = 72,766 bytes
- 35 pages + 1,086 bytes
- Internal fragmentation of  $2,048 - 1,086 = 962$  bytes
- Worst case fragmentation = 1 frame – 1 byte
- On average fragmentation =  $1 / 2$  frame size
- So small frame sizes desirable?
- But each page table entry takes memory to track
- Page sizes growing over time
  - Solaris supports two page sizes – 8 KB and 4 MB

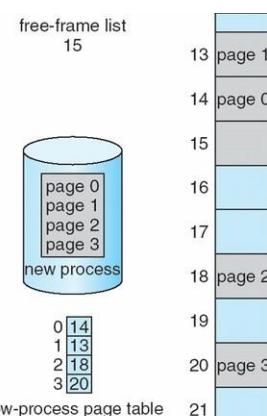


## Free Frames (Self Study)



(a)

Before allocation



(b)

After allocation





## Implementation of Page Table (Self Study)

- 
- Page table is kept in main memory
    - **Page-table base register (PTBR)** points to the page table
    - **Page-table length register (PTLR)** indicates size of the page table
  - In this scheme every data/instruction access requires two memory accesses
    - One for the page table and one for the data / instruction
  - The two-memory access problem can be solved by the use of a special fast-lookup hardware cache called **translation look-aside buffers (TLBs)** (also called **associative memory**).

Operating System Concepts – 10<sup>th</sup> Edition

9.223

Silberschatz, Galvin and Gagne ©2018

223



## Translation Look-Aside Buffer (Self Study)

- 
- Some TLBs store **address-space identifiers (ASIDs)** in each TLB entry – uniquely identifies each process to provide address-space protection for that process
    - Otherwise need to flush at every context switch
  - TLBs typically small (64 to 1,024 entries)
  - On a TLB miss, value is loaded into the TLB for faster access next time
    - Replacement policies must be considered
    - Some entries can be **wired down** for permanent fast access

Operating System Concepts – 10<sup>th</sup> Edition

9.224

Silberschatz, Galvin and Gagne ©2018

224



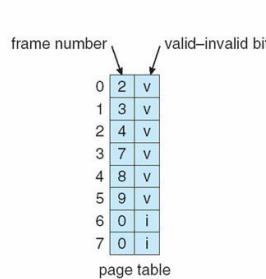
## Memory Protection (Self Study)

- Memory protection implemented by associating protection bit with each frame to indicate if read-only or read-write access is allowed
  - Can also add more bits to indicate page execute-only, and so on
- **Valid-invalid** bit attached to each entry in the page table:
  - “valid” indicates that the associated page is in the process’ logical address space, and is thus a legal page
  - “invalid” indicates that the page is not in the process’ logical address space
  - Or use **page-table length register (PTLR)**
- Any violations result in a trap to the kernel



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

|        |        |
|--------|--------|
| 00000  | page 0 |
|        | page 1 |
|        | page 2 |
|        | page 3 |
|        | page 4 |
| 10,468 | page 5 |
| 12,287 |        |



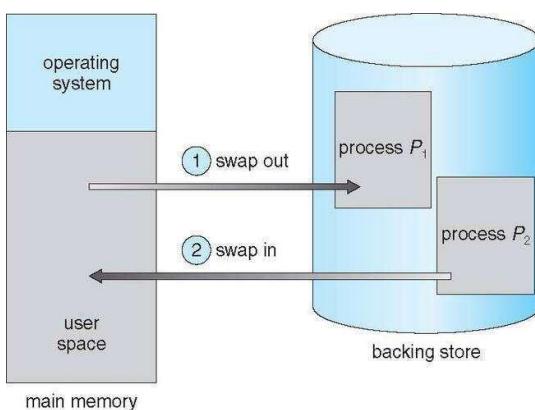


## Swapping

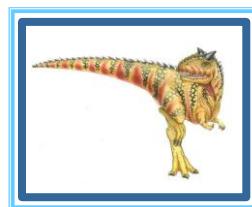
- A process can be **swapped** temporarily out of memory to a backing store, and then brought **back** into memory for continued execution
  - Total physical memory space of processes can exceed physical memory
- **Backing store** – fast disk large enough to accommodate copies of all memory images for all users; must provide direct access to these memory images
- **Roll out, roll in** – swapping variant used for priority-based scheduling algorithms; lower-priority process is swapped out so higher-priority process can be loaded and executed
- Major part of swap time is transfer time; total transfer time is directly proportional to the amount of memory swapped
- System maintains a **ready queue** of ready-to-run processes which have memory images on disk



## Schematic View of Swapping



## Chapter 5 (Cont'd): Virtual Memory



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

229



## Chapter 5 (Cont'd): Virtual Memory

- Background
- Demand Paging
- Copy-on-Write
- Page Replacement
- Allocation of Frames
- Thrashing

Operating System Concepts – 10<sup>th</sup> Edition

9.230

Silberschatz, Galvin and Gagne ©2018



230

115



## 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

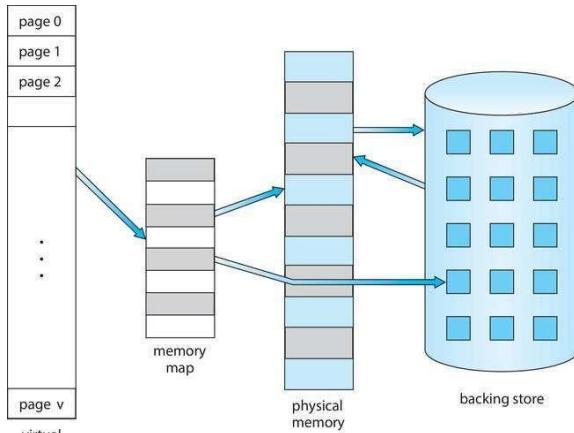


## Virtual memory (Cont.)

- **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



## Virtual Memory That is Larger Than Physical Memory (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

9.233

Silberschatz, Galvin and Gagne ©2018

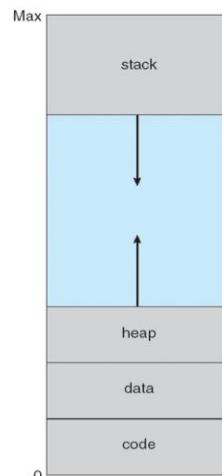


233



## Virtual-address Space (Self Study)

- Usually design logical address space for stack to start at Max logical address and grow “down” while heap grows “up”
  - Maximizes address space use
  - Unused address space between the two is hole
    - No physical memory needed until heap or stack grows to a given new page
- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc.
- System libraries shared via mapping into virtual address space
- Shared memory by mapping pages read-write into virtual address space
- Pages can be shared during `fork()`, speeding process creation



Operating System Concepts – 10<sup>th</sup> Edition

9.234

Silberschatz, Galvin and Gagne ©2018

234



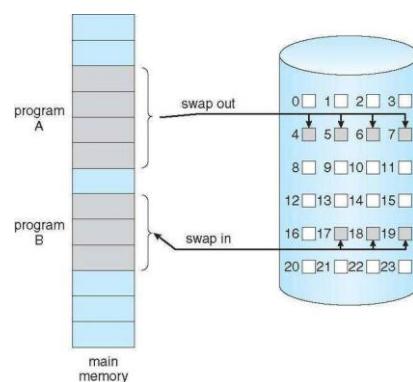
## Demand Paging

- 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**



## 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)





## Basic Concepts (Self Study)

- With swapping, pager guesses which pages will be used before swapping out again
- Instead, pager brings in only those pages into memory
- How to determine that set of pages?
  - Need new MMU functionality to implement demand paging
- If pages needed are already **memory resident**
  - No difference from non demand-paging
- If page needed and not memory resident
  - Need to detect and load the page into memory from storage
    - Without changing program behavior
    - Without programmer needing to change code



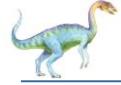
## Valid-Invalid Bit

- 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:

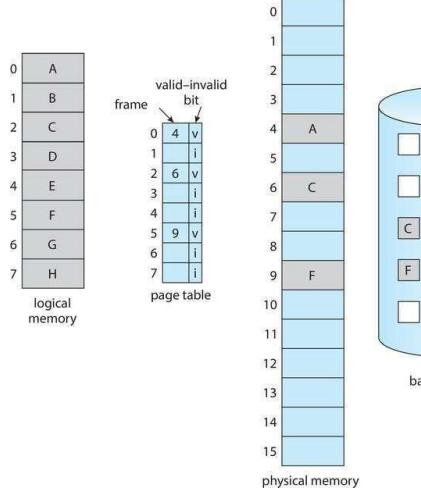
| Frame # | valid-invalid bit |
|---------|-------------------|
|         |                   |
|         | v                 |
|         | v                 |
|         | v                 |
| ...     | i                 |
|         | i                 |
|         | i                 |

- 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



Operating System Concepts – 10<sup>th</sup> Edition

9.239

Silberschatz, Galvin and Gagne ©2018



239



## Steps in Handling Page Fault (Self Study)

1. If there is a reference to a page, first reference to that page will trap to operating system
  - Page fault
2. Operating system looks at another table to decide:
  - Invalid reference  $\Rightarrow$  abort
  - Just not in memory
3. Find free frame
4. Swap page into frame via scheduled disk operation
5. Reset tables to indicate page now in memory  
Set validation bit = **v**
6. Restart the instruction that caused the page fault

Operating System Concepts – 10<sup>th</sup> Edition

9.240

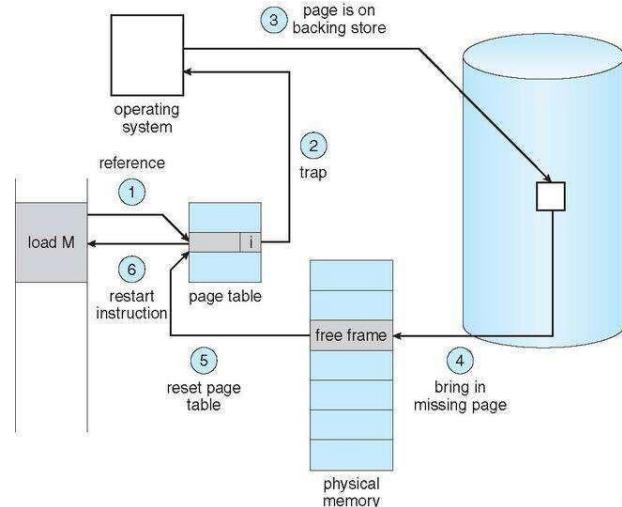
Silberschatz, Galvin and Gagne ©2018



240

120

## Steps in Handling a Page Fault (Cont.) (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

9.241

Silberschatz, Galvin and Gagne ©2018

241

## Aspects of Demand Paging (Self Study)

- Extreme case – start process with *no* pages in memory
  - OS sets instruction pointer to first instruction of process, non-memory-resident -> page fault
  - And for every other process pages on first access
  - Pure demand paging**
- Actually, a given instruction could access multiple pages -> multiple page faults
  - Consider fetch and decode of instruction which adds 2 numbers from memory and stores result back to memory
  - Pain decreased because of **locality of reference**
- Hardware support needed for demand paging
  - Page table with valid / invalid bit
  - Secondary memory (swap device with **swap space**)
  - Instruction restart

Operating System Concepts – 10<sup>th</sup> Edition

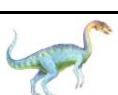
9.242

Silberschatz, Galvin and Gagne ©2018

242

## Stages in Demand Paging – Worse Case (Self Study)

1. Trap to the operating system
2. Save the user registers and process state
3. Determine that the interrupt was a page fault
4. Check that the page reference was legal and determine the location of the page on the disk
5. Issue a read from the disk to a free frame:
  - a) Wait in a queue for this device until the read request is serviced
  - b) Wait for the device seek and/or latency time
  - c) Begin the transfer of the page to a free frame



## Stages in Demand Paging (Cont.) (Self Study)

6. While waiting, allocate the CPU to some other user
7. Receive an interrupt from the disk I/O subsystem (I/O completed)
8. Save the registers and process state for the other user
9. Determine that the interrupt was from the disk
10. Correct the page table and other tables to show page is now in memory
11. Wait for the CPU to be allocated to this process again
12. Restore the user registers, process state, and new page table, and then resume the interrupted instruction



## Performance of Demand Paging (Self Study)



- Three major activities
  - Service the interrupt – careful coding means just several hundred instructions needed
  - Read the page – lots of time
  - Restart the process – again just a small amount of time
- Page Fault Rate  $0 \leq p \leq 1$ 
  - if  $p = 0$  no page faults
  - if  $p = 1$ , every reference is a fault
- Effective Access Time (EAT)
$$\begin{aligned} EAT &= (1 - p) \times \text{memory access} \\ &\quad + p (\text{page fault overhead} \\ &\quad + \text{swap page out} \\ &\quad + \text{swap page in}) \end{aligned}$$



## Demand Paging Example (Self Study)

- Memory access time = 200 nanoseconds
- Average page-fault service time = 8 milliseconds
- $$\begin{aligned} EAT &= (1 - p) \times 200 + p (8 \text{ milliseconds}) \\ &= (1 - p \times 200 + p \times 8,000,000 \\ &= 200 + p \times 7,999,800 \end{aligned}$$
- If one access out of 1,000 causes a page fault, then
$$EAT = 8.2 \text{ microseconds.}$$
This is a slowdown by a factor of 40!!
- If want performance degradation < 10 percent
  - $220 > 200 + 7,999,800 \times p$   
 $20 > 7,999,800 \times p$
  - $p < .0000025$
  - < one page fault in every 400,000 memory accesses





## 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

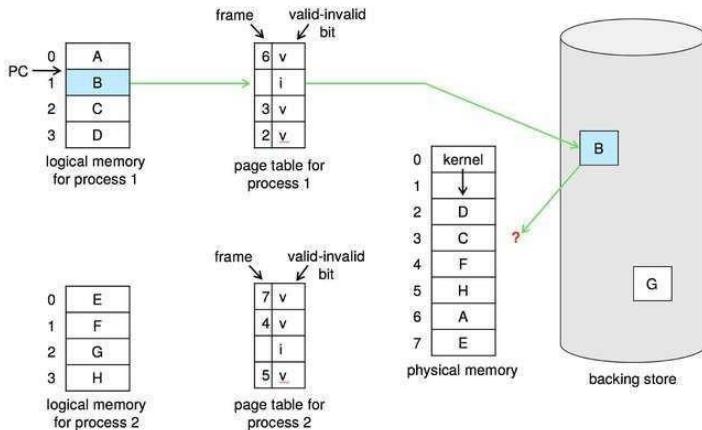


## Page Replacement (Self Study)

- Prevent **over-allocation** of memory by modifying page-fault service routine to include page replacement
- Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk
- Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory



## Need For Page Replacement (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

9.249

Silberschatz, Galvin and Gagne ©2018

249

## 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

Operating System Concepts – 10<sup>th</sup> Edition

9.250

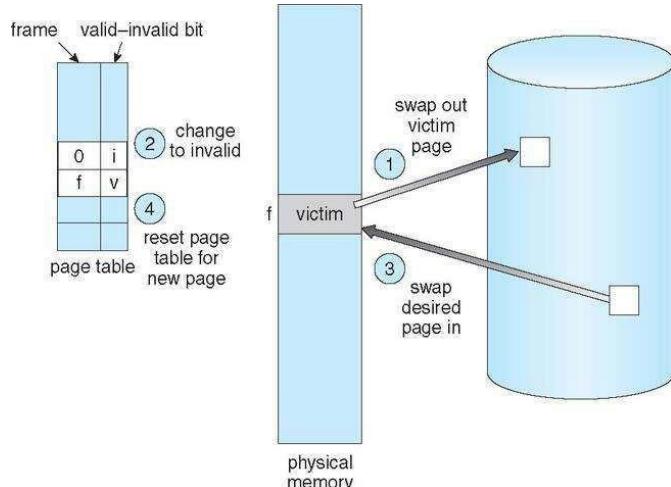
Silberschatz, Galvin and Gagne ©2018

250

125



## Page Replacement (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

9.251

Silberschatz, Galvin and Gagne ©2018



251



## 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**

Operating System Concepts – 10<sup>th</sup> Edition

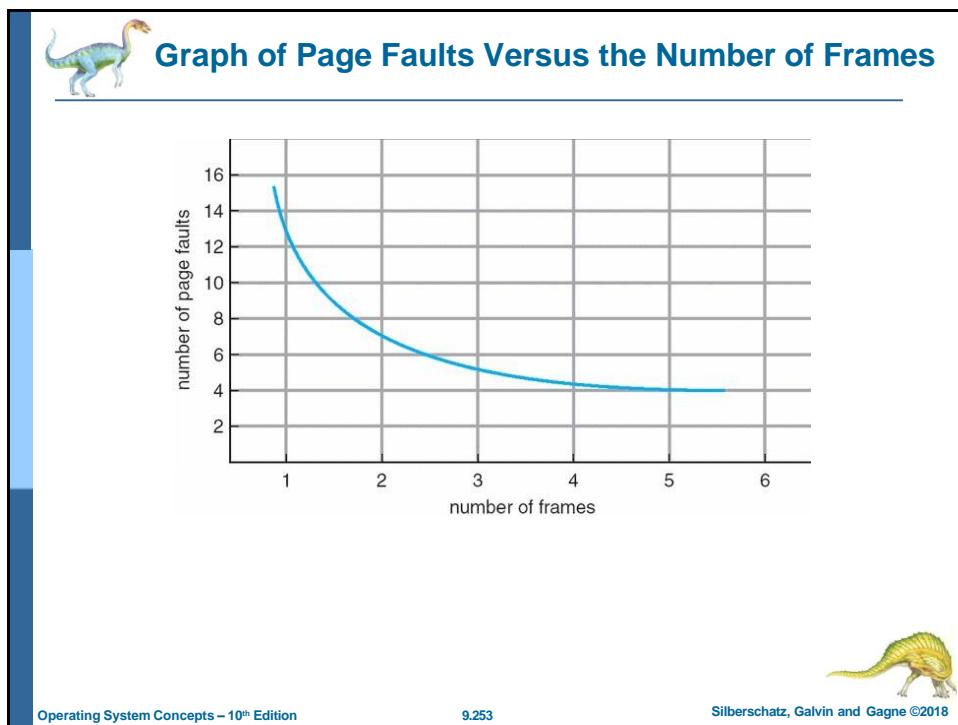
9.252

Silberschatz, Galvin and Gagne ©2018

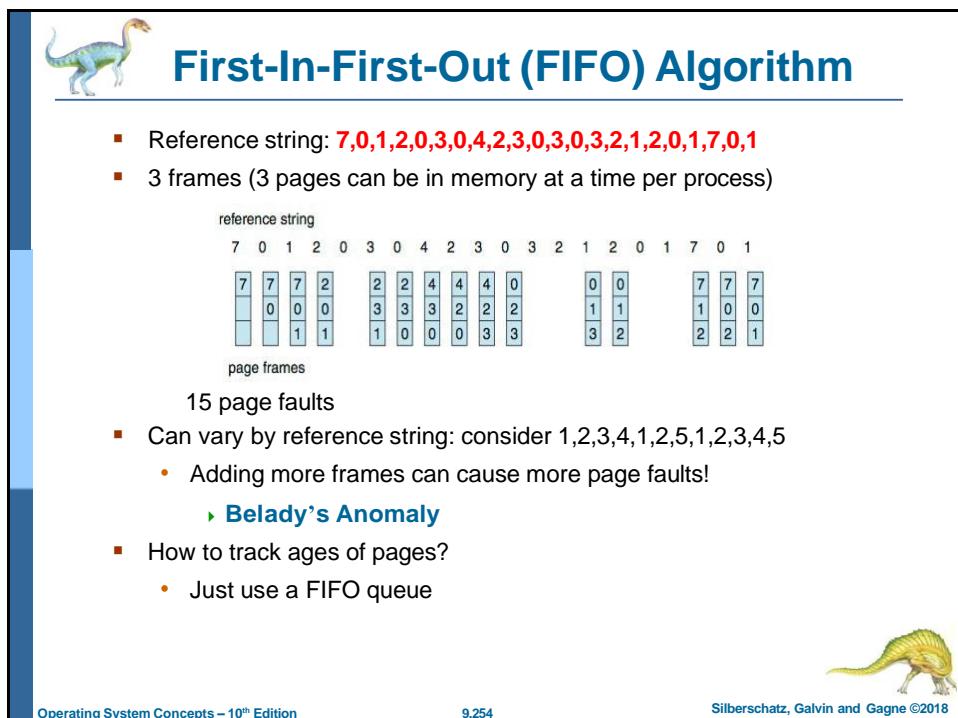


252

126



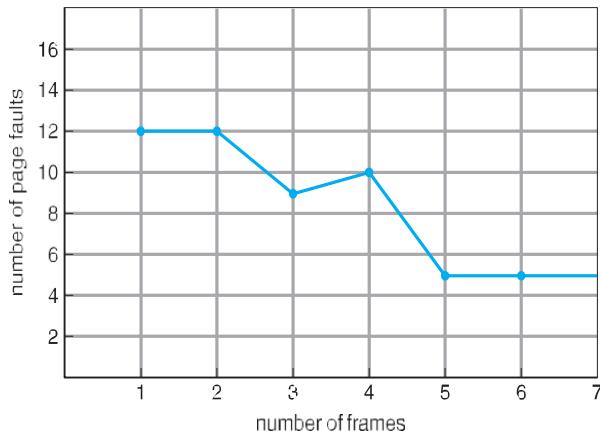
253



254



## FIFO Illustrating Belady's Anomaly



Operating System Concepts – 10<sup>th</sup> Edition

9.255

Silberschatz, Galvin and Gagne ©2018

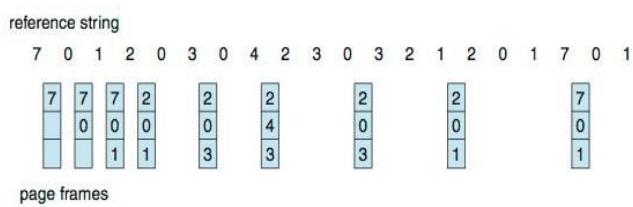


255



## 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 – 10<sup>th</sup> Edition

9.256

Silberschatz, Galvin and Gagne ©2018

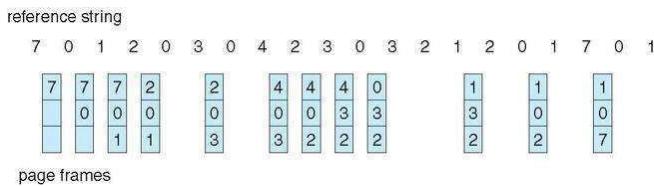


256

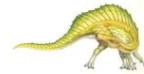


## Least Recently Used (LRU) Algorithm

- 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



- 12 faults – better than FIFO but worse than OPT
- Generally good algorithm and frequently used
- But how to implement?



## LRU Algorithm (Cont.)

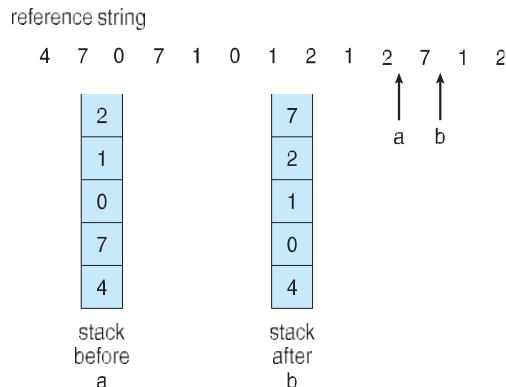
- Counter implementation
  - Every page entry has a counter; every time page is referenced through this entry, copy the clock into the counter
  - When a page needs to be changed, look at the counters to find smallest value
    - ▶ Search through table needed
- Stack implementation
  - Keep a stack of 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





## 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



## Thrashing

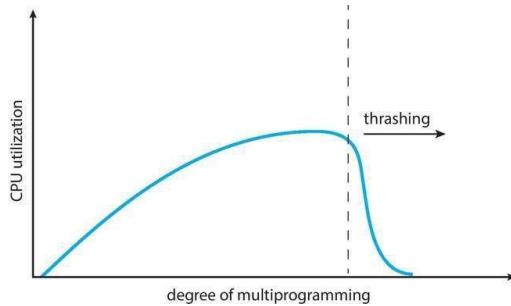
- If a process does not have “enough” pages, the page-fault rate is very high
  - Page fault to get page
  - Replace existing frame
  - But quickly need replaced frame back
  - This leads to:
    - ▶ Low CPU utilization
    - ▶ Operating system thinking that it needs to increase the degree of multiprogramming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



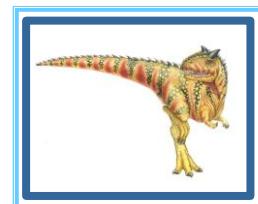


## Thrashing (Cont.)

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



## Chapter 6: File Management





## Outline

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



## 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



Operating System Concepts – 10<sup>th</sup> Edition

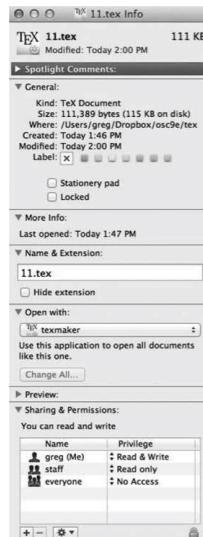
13.265

Silberschatz, Galvin and Gagne ©2018

265



## (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

13.266

Silberschatz, Galvin and Gagne ©2018

266



## 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



## Open Files (Self Study)

- Several pieces of data are needed to manage open files:
  - **Open-file table**: tracks open files
  - File pointer: pointer to last read/write location, per process that has the file open
  - **File-open count**: counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it
  - Disk location of the file: cache of data access information
  - Access rights: per-process access mode information





## 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



## (Self Study)

| file type      | usual extension          | function  |
|----------------|--------------------------|---|
| executable     | exe, com, bin or none    | ready-to-run machine-language program   |
| object         | obj, o                   | compiled, machine language, not linked  |
| source code    | c, cc, java, pas, asm, a | source code in various languages  |
| batch          | bat, sh                  | commands to the command interpreter   |
| text           | txt, doc                 | textual data, documents   |
| word processor | wp, tex, rtf, doc        | various word-processor formats  |
| library        | lib, a, so, dll          | libraries of routines for programmers   |
| print or view  | ps, pdf, jpg             | ASCII or binary file in a format for printing or viewing                            |
| archive        | arc, zip, tar            | related files grouped into one file, sometimes compressed, for archiving or storage |
| multimedia     | mpeg, mov, rm, mp3, avi  | binary file containing audio or A/V information                                     |





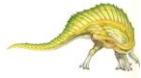
## File Structure

- 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

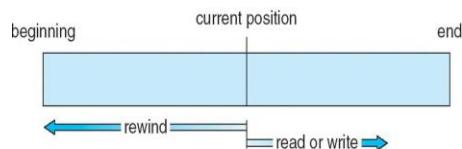
- A file is fixed length **logical records**
- **Sequential Access**
- **Direct Access**
- **Other Access Methods**





## Sequential Access

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



Operating System Concepts – 10<sup>th</sup> Edition

13.273

Silberschatz, Galvin and Gagne ©2018



273



## Direct Access

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

Operating System Concepts – 10<sup>th</sup> Edition

13.274

Silberschatz, Galvin and Gagne ©2018



274




## Simulation of Sequential Access on Direct-access File (Self Study)

| sequential access | implementation for direct access   |
|-------------------|------------------------------------|
| <i>reset</i>      | $cp = 0;$                          |
| <i>read next</i>  | <i>read cp;</i><br>$cp = cp + 1;$  |
| <i>write next</i> | <i>write cp;</i><br>$cp = cp + 1;$ |




Operating System Concepts – 10<sup>th</sup> Edition

13.275

Silberschatz, Galvin and Gagne ©2018

275




## (Self Study)

- Can be other access methods built on top of base methods
- General involve creation of an **index** for the file
- Keep index in memory for fast determination of location of data to be operated on (consider Universal Produce Code (UPC code) plus record of data about that item)
- If the index is too large, create an in-memory index, which is an index of a disk index
- IBM indexed sequential-access method (ISAM)
  - Small master index, points to disk blocks of secondary index
  - File kept sorted on a defined key
  - All done by the OS
- VMS operating system provides index and relative files as another example (see next slide)




Operating System Concepts – 10<sup>th</sup> Edition

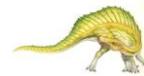
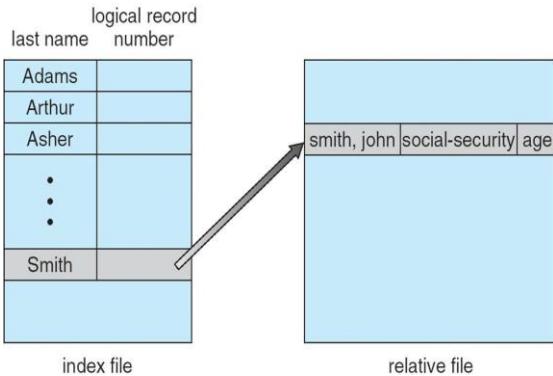
13.276

Silberschatz, Galvin and Gagne ©2018

276



## (Self Study)



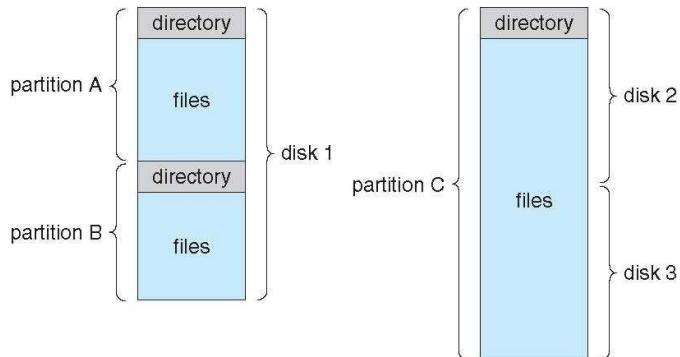
## Disk Structure

- 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



Operating System Concepts – 10<sup>th</sup> Edition

13.279

Silberschatz, Galvin and Gagne ©2018



279



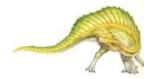
## (Self Study)

- We mostly talk of general-purpose file systems
- But systems frequently have many file systems, some general- and some special- purpose
- Consider Solaris has
  - tmpfs – memory-based volatile FS for fast, temporary I/O
  - objfs – interface into kernel memory to get kernel symbols for debugging
  - ctfs – contract file system for managing daemons
  - lofs – loopback file system allows one FS to be accessed in place of another
  - procfs – kernel interface to process structures
  - ufs, zfs – general purpose file systems

Operating System Concepts – 10<sup>th</sup> Edition

13.280

Silberschatz, Galvin and Gagne ©2018



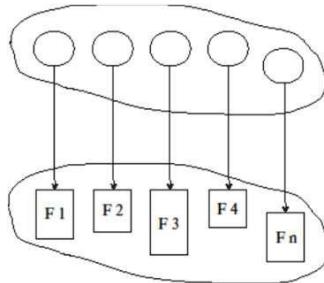
280

140



## 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

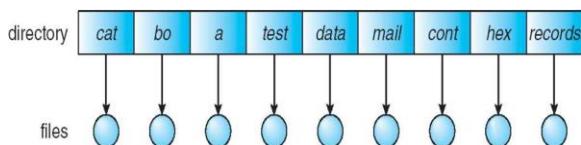
The directory is organized logically to obtain

- 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



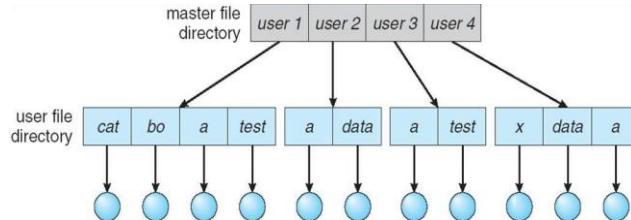
- Naming problem
- Grouping problem





## Two-Level Directory

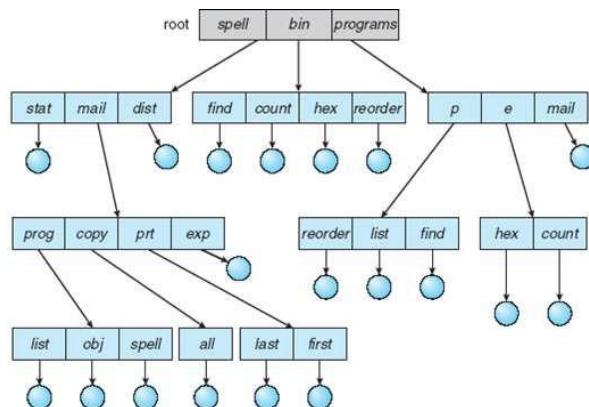
- Separate directory for each user



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



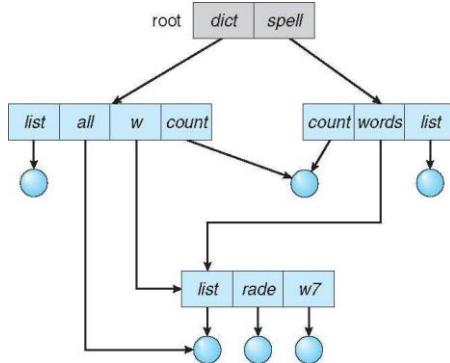
## Tree-Structured Directories (Self Study)





## Acyclic-Graph Directories (Self Study)

- Have shared subdirectories and files
- Example



Operating System Concepts – 10<sup>th</sup> Edition

13.287

Silberschatz, Galvin and Gagne ©2018



287



## (Self Study)

- Two different names (aliasing)
  - If **dict** deletes **w/list** ⇒ dangling pointer
- Solutions:
- Backpointers, so we can delete all pointers.
    - Variable size records a problem
  - Backpointers using a daisy chain organization
  - Entry-hold-count solution
- New directory entry type
    - **Link** – another name (pointer) to an existing file
    - **Resolve the link** – follow pointer to locate the file

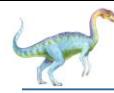
Operating System Concepts – 10<sup>th</sup> Edition

13.288

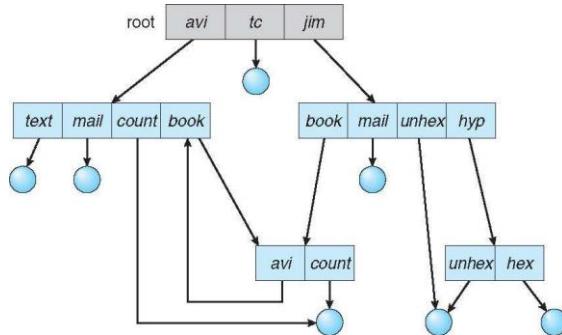
Silberschatz, Galvin and Gagne ©2018



288



## General Graph Directory (Self Study)

Operating System Concepts – 10<sup>th</sup> Edition

13.289

Silberschatz, Galvin and Gagne ©2018

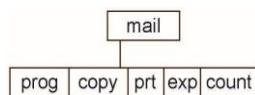


289



## 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”

Operating System Concepts – 10<sup>th</sup> Edition

13.290

Silberschatz, Galvin and Gagne ©2018

290



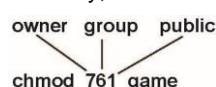
## Protection

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



## Study)

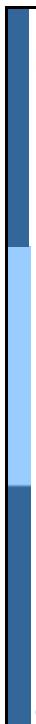
- Mode of access: read, write, execute
- Three classes of users on Unix / Linux
  - a) owner access      7       $\Rightarrow$  1 1      1  
                                  RWX
  - b) group access      6       $\Rightarrow$  1 1 0  
                                  RWX
  - c) public access      1       $\Rightarrow$  0 0 1
- Ask manager to create a group (unique name), say G, and add some users to the group.
- For a file (say *game*) or subdirectory, define an appropriate access.



- Attach a group to a file

`chgrp      G      game`





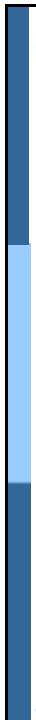
**(Self Study)**

| File Type  | Owner | Group   | Last Modified | Last Accessed | File Name     |
|------------|-------|---------|---------------|---------------|---------------|
| -rw-rw-r-- | 1 pbg | staff   | 31200         | Sep 3 08:30   | intro.ps      |
| drwx-----  | 5 pbg | staff   | 512           | Jul 8 09:33   | private/      |
| drwxrwxr-x | 2 pbg | staff   | 512           | Jul 8 09:35   | doc/          |
| drwxrwx--- | 2 pbg | student | 512           | Aug 3 14:13   | student-proj/ |
| -rw-r--r-- | 1 pbg | staff   | 9423          | Feb 24 2003   | program.c     |
| -rwxr-xr-x | 1 pbg | staff   | 20471         | Feb 24 2003   | program       |
| drwx--x--x | 4 pbg | faculty | 512           | Jul 31 10:31  | lib/          |
| drwx-----  | 3 pbg | staff   | 1024          | Aug 29 06:52  | mail/         |
| drwxrwxrwx | 3 pbg | staff   | 512           | Jul 8 09:35   | test/         |

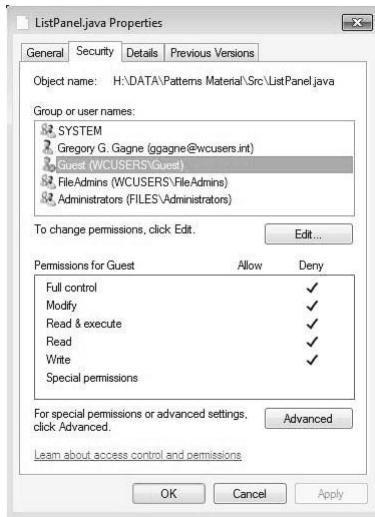


Operating System Concepts – 10<sup>th</sup> Edition      13.293      Silberschatz, Galvin and Gagne ©2018

293



**Windows 7 Access-Control List Management  
(Self Study)**



The dialog box shows the security properties for the file 'ListPanel.java'. The 'Guest' user has full control (Allow) for all permissions listed.

| Permission          | Allow |
|---------------------|-------|
| Full control        | ✓     |
| Modify              | ✓     |
| Read & execute      | ✓     |
| Read                | ✓     |
| Write               | ✓     |
| Special permissions | ✓     |

For special permissions or advanced settings, click Advanced.

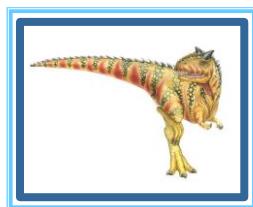
OK Cancel Apply



Operating System Concepts – 10<sup>th</sup> Edition      13.294      Silberschatz, Galvin and Gagne ©2018

294

# Chapter 6 (Cont'd): File System Implementation



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

295



## Outline

- File-System Structure
- File-System Operations
- Allocation Methods

Operating System Concepts – 10<sup>th</sup> Edition

13.296

Silberschatz, Galvin and Gagne ©2018

296



148

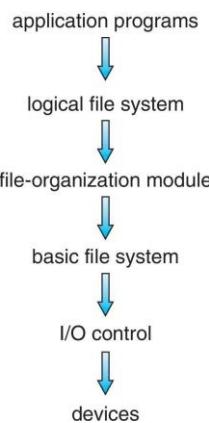


## File-System Structure

- 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



Operating System Concepts – 10<sup>th</sup> Edition

13.299

Silberschatz, Galvin and Gagne ©2018

299



## 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



Operating System Concepts – 10<sup>th</sup> Edition

13.300

Silberschatz, Galvin and Gagne ©2018

300

150



## (Self Study)

- Many file systems, sometimes many within an operating system
  - Each with its own format:
  - CD-ROM is ISO 9660;
  - Unix has **UFS**, FFS;
  - Windows has FAT, FAT32, NTFS as well as floppy, CD, DVD Blu-ray,
  - Linux has more than 130 types, with **extended file system** ext3 and ext4 leading; plus distributed file systems, etc.)
  - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE



## File-System Operations (Self Study)

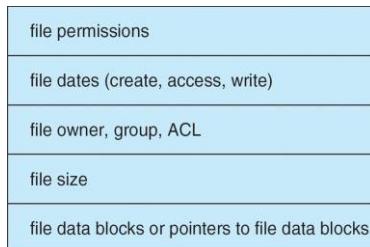
- We have system calls at the API level, but how do we implement their functions?
  - On-disk and in-memory structures
- **Boot control block** contains info needed by system to boot OS from that volume
  - Needed if volume contains OS, usually first block of volume
- **Volume control block (superblock, master file table)** contains volume details
  - Total # of blocks, # of free blocks, block size, free block pointers or array
- Directory structure organizes the files
  - Names and inode numbers, master file table





## File Control Block (FCB) (Self Study)

- OS maintains **FCB** per file, which contains many details about the file
  - Typically, inode number, permissions, size, dates
  - Example



Operating System Concepts – 10<sup>th</sup> Edition

13.303

Silberschatz, Galvin and Gagne ©2018



303



## Allocation Method

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

Operating System Concepts – 10<sup>th</sup> Edition

13.304

Silberschatz, Galvin and Gagne ©2018



304

152



## 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**



Operating System Concepts – 10<sup>th</sup> Edition

13.305

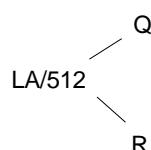
Silberschatz, Galvin and Gagne ©2018

305



## Contiguous Allocation (Cont.)

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



| directory |       |        |
|-----------|-------|--------|
| file      | start | length |
| count     | 0     | 2      |
| tr        | 14    | 3      |
| mail      | 19    | 6      |
| list      | 28    | 4      |
| f         | 6     | 2      |

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



Operating System Concepts – 10<sup>th</sup> Edition

13.306

Silberschatz, Galvin and Gagne ©2018

306



## 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 – 10<sup>th</sup> Edition

13.307

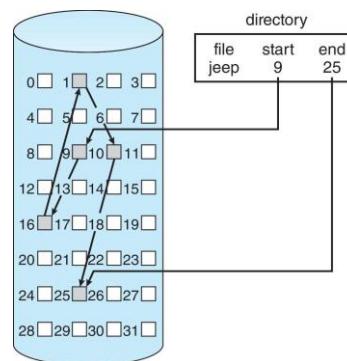
Silberschatz, Galvin and Gagne ©2018

307



## Linked Allocation Example

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



Operating System Concepts – 10<sup>th</sup> Edition

13.308

Silberschatz, Galvin and Gagne ©2018

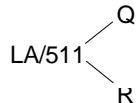
308

154



## Linked Allocation (Cont.)

- Mapping



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



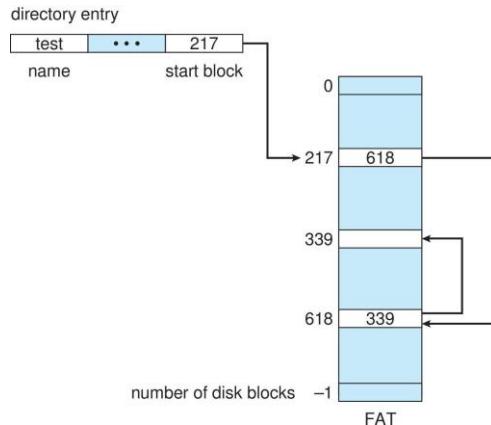
## FAT Allocation Method

- 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



Operating System Concepts – 10<sup>th</sup> Edition

13.311

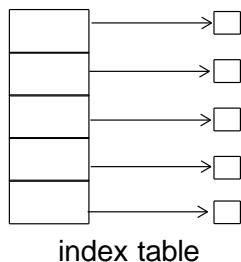
Silberschatz, Galvin and Gagne ©2018

311



## Indexed Allocation Method

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



Operating System Concepts – 10<sup>th</sup> Edition

13.312

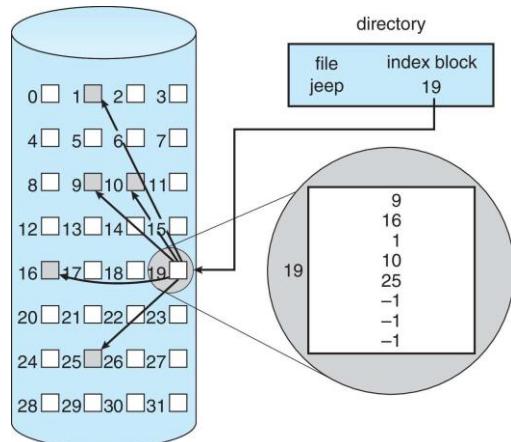
Silberschatz, Galvin and Gagne ©2018

312

156



## Example of Indexed Allocation



Operating System Concepts – 10<sup>th</sup> Edition

13.313

Silberschatz, Galvin and Gagne ©2018

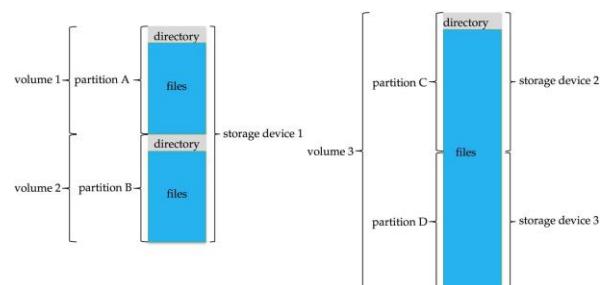


313



## File System (Self Study)

- General-purpose computers can have multiple storage devices
- Devices can be sliced into partitions, which hold volumes
- Volumes can span multiple partitions
- Each volume usually formatted into a file system
- # of file systems varies, typically dozens available to choose from
- Typical storage device organization:



Operating System Concepts – 10<sup>th</sup> Edition

13.314

Silberschatz, Galvin and Gagne ©2018



314

157



## (Self Study)

|                   |       |
|-------------------|-------|
| /                 | ufs   |
| /devices          | devfs |
| /dev              | dev   |
| /system/contract  | ctfs  |
| /proc             | proc  |
| /etc/mnttab       | mntfs |
| /etc/svc/volatile | tmpfs |
| /system/object    | objfs |
| /lib/libc.so.1    | lofs  |
| /dev/fd           | fd    |
| /var              | ufs   |
| /tmp              | tmpfs |
| /var/run          | tmpfs |
| /opt              | ufs   |
| /zpbge            | zfs   |
| /zpbge/backup     | zfs   |
| /export/home      | zfs   |
| /var/mail         | zfs   |
| /var/spool/mqueue | zfs   |
| /zpbg             | zfs   |
| /zpbg/zones       | zfs   |



## File Sharing

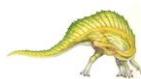
- Allows multiple users / systems access to the same files
- Permissions / protection must be implemented and accurate
  - Most systems provide concepts of owner, group member
  - Must have a way to apply these between systems





## (Self Study)

- **Virtual File Systems (VFS)** on Unix provide an object-oriented way of implementing file systems
- VFS allows the same system call interface (the API) to be used for different types of file systems
  - Separates file-system generic operations from implementation details
  - Implementation can be one of many file systems types, or network file system
    - Implements **vnodes** which hold inodes or network file details
  - Then dispatches operation to appropriate file system implementation routines



Operating System Concepts – 10<sup>th</sup> Edition

13.317

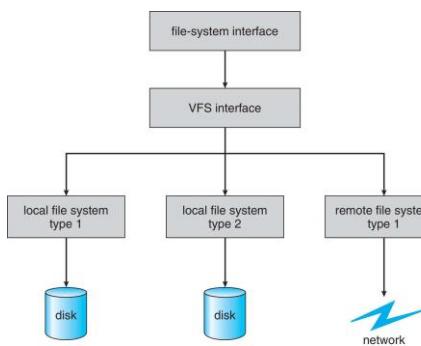
Silberschatz, Galvin and Gagne ©2018

317



## (Self Study)

- The API is to the VFS interface, rather than any specific type of file system
- Example



Operating System Concepts – 10<sup>th</sup> Edition

13.318

Silberschatz, Galvin and Gagne ©2018

318

159



## (Self Study)

- For example, Linux has four object types:
  - **inode, file, superblock, dentry**
- VFS defines set of operations on the objects that must be implemented
  - Every object has a pointer to a function table
    - ▶ Function table has addresses of routines to implement that function on that object
    - ▶ For example:
      - ▶ • `int open(...)`—Open a file
      - ▶ • `int close(...)`—Close an already-open file
      - ▶ • `ssize_t read(...)`—Read from a file
      - ▶ • `ssize_t write(...)`—Write to a file
      - ▶ • `int mmap(...)`—Memory-map a file



## Remote File Systems (Self Study)

- Sharing of files across a network
- First method involved manually sharing each file – programs like **ftp**
- Second method uses a **distributed file system (DFS)**
  - Remote directories visible from local machine
- Third method – **World Wide Web**
  - A bit of a revision to first method
  - Use browser to locate file/files and download /upload
  - **Anonymous** access doesn't require authentication





## (Self Study)

- An implementation and a specification of a software system for accessing remote files across LANs (or WANs)
- The implementation originally part of SunOS operating system, now industry standard / very common
- Can use unreliable datagram protocol (UDP/IP) or TCP/IP, over Ethernet or other networks



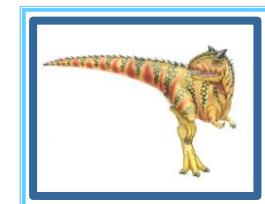
Operating System Concepts – 10<sup>th</sup> Edition

13.321

Silberschatz, Galvin and Gagne ©2018

321

## Chapter 6 (Cont'd): OS Security



Operating System Concepts – 10<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2018

322

161



## 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



## 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





## 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



## 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?





## Program Threats (Self Study)

- Many variations, many names
- **Trojan Horse**
  - Code segment that misuses its environment
  - Exploits mechanisms for allowing programs written by users to be executed by other users
  - **Spyware, pop-up browser windows, covert channels**
  - Up to 80% of spam delivered by spyware-infected systems
- **Trap Door**
  - Specific user identifier or password that circumvents normal security procedures
  - Could be included in a compiler
  - How to detect them?



Operating System Concepts – 10<sup>th</sup> Edition

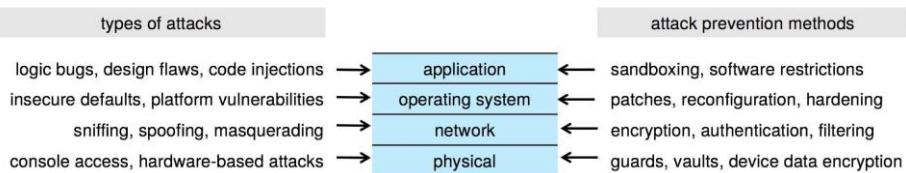
13.327

Silberschatz, Galvin and Gagne ©2018

327



## Four-layered Model of Security (Self Study)



Operating System Concepts – 10<sup>th</sup> Edition

13.328

Silberschatz, Galvin and Gagne ©2018

328

164



## Program Threats (Cont.) (Self Study)

- **Malware** – Software designed to exploit, disable, or damage computer
- **Trojan Horse** – Program that acts in a clandestine manner
  - **Spyware** – Program frequently installed with legitimate software to display adds, capture user data
  - **Ransomware** – locks up data via encryption, demanding payment to unlock it
- Others include trap doors, logic bombs
- All try to violate the Principle of Least Privilege

### THE PRINCIPLE OF LEAST PRIVILEGE

"The principle of least privilege. Every program and every privileged user of the system should operate using the least amount of privilege necessary to complete the job. The purpose of this principle is to reduce the number of potential interactions among privileged programs to the minimum necessary to operate correctly, so that one may develop confidence that unintentional, unwanted, or improper uses of privilege do not occur."—Jerome H. Saltzer, describing a design principle of the Multics operating system in 1974: <https://pdfs.semanticscholar.org/1c8d/06510ad449ad24fbdd164f8008cc730cab47.pdf>.

- ... access Tool (RAT) for repeated access



## End of Chapter 6

