

Daniel **Fitz**
(43961229)



University Of Queensland
COMP3301 – Operating Systems Architecture

COMP3301 Lecture Notes



Table of Contents

Course Outline	9
Introduction	9
What is an operating system?	9
What Operating Systems Do	9
Operating System Definition	9
Common Functions of Interrupts	10
Interrupt Handling	10
Computer-System Architecture	10
Operating System Structure	10
Operating-System Operations	10
Transition from User to Kernel Mode	11
Process Management	11
Memory Management	11
Storage Management	11
IO Subsystem	11
Protection and Security	12
Open-Source Operating Systems	12
Operating-System Structures	12
Operating System Services	12
User Operating System Interface	12
CLI	13
GUI	13
Touchscreen	13
System Calls	13
System Call Parameter Passing	13
Types of System Calls	13
System Programs	14
Operating System Design and Implementation	15
Implementation	15
UNIX	15
Layered Approach	15
Microkernel System Structure	16
Modules	16
System Boot	16
Processes	16
Process Concept	16
Process State	17
Process Control Block (PCB)	17
Threads	17
Process Scheduling	17

Schedulers	18
Addition of Medium Term Scheduling	18
Multitasking in Mobile Systems	18
Context Switch	18
Operations on Processes	18
Process Creation	18
Process Termination	19
Interprocess Communication	19
Cooperating Processes	19
Producer-Consumer Problem	20
Interprocess Communication - Message Passing	20
Direct Communication	20
Indirect Communication	20
Synchronization	21
Buffer	21
Examples of IPC Systems	21
POSIX	21
Communication in Client-Server Systems	21
Sockets	21
Remote Procedure Calls	21
Pipes	22
Threads	22
Overview	22
Motivation	22
Benefits	23
Multicore Programming	23
Amdahls Law	23
User Threads and Kernel Threads	24
Multithreading Models	24
Many-to-One	24
One-to-One	24
Many-to-Many	24
Thread Libraries	24
Pthreads	24
Java Threads	25
Implicit Threading	25
Thread Pools	25
OpenMP	25
Grand Central Dispatch	25
Threading Issues	26
Semantics of fork() and exec()	26

Signal Handling	26
Thread-Local Storage	26
Linux Threads	26
CPU Scheduling	27
Basic Concepts	27
CPU Scheduler	27
Dispatcher	28
Scheduling Criteria	28
Scheduling Algorithms	28
Criteria	28
First-Come, First-Served (FCFS) Scheduling	28
Shortest-Job-First (SJF) Scheduling	29
Example of Shortest-remaining-time-first	29
Priority Scheduling	29
Example of Priority Scheduling	29
Round Robin	30
Multilevel Queue	30
Multilevel Feedback Queue	30
Thread Scheduling	30
Multiple-Processor Scheduling	31
Load Balancing	31
Multicore Processors	31
Virtualization and Scheduling	31
Real-Time CPU Scheduling	31
Priority-based Scheduling	31
Rate Monotonic Scheduling	32
Earliest Deadline First Scheduling (EDF)	32
Proportional Share Scheduling	32
Operating Systems Examples	32
Linux Scheduling	32
Windows Scheduling	32
Solaris Scheduling	33
Algorithm Evaluation	33
Deterministic Evaluation	33
Queueing Models	33
Little's Law	34
Simulations	34
Implementation	34
Process Synchronization	34
Background	34
The Critical-Section Problem	34

Critical Section	35
Solution to Critical-Section Problem	35
Petersons Solution	35
Algorithm for Process P_i	35
Synchronization Hardware	36
Mutex Locks	36
Semaphore	36
Semaphore Usage	36
Semaphore Implementation	36
Semaphore Implementation with no Busy waiting	37
Deadlock and Starvation	37
Classic Problems of Synchronization	37
Bounded-Buffer Problem	37
Readers-Writers Problem	38
Dining-Philosophers Problem	38
Problems with Semaphores	38
Monitors	38
Condition Variables	39
Condition Variables Choices	39
Linux Synchronisation	39
Atomic Transactions	39
Deadlocks	40
System Model	40
Deadlock Characterization	40
Deadlock with Mutex Locks	40
Resource-Allocation Graph	40
Methods for Handling Deadlocks	41
Deadlock Prevention	41
Deadlock Avoidance	41
Safe State	42
Avoidance Algorithms	42
Resource-Allocation Graph Scheme	42
Bankers Algorithm	42
Deadlock Detection	42
Single Instance of Each Resource Type	42
Detection-Algorithm Usage	43
Recovery from Deadlock	43
Main Memory	43
Background	43
Base and Limit Registers	43
Address Binding	44

Logical vs. Physical Address Space	44
Memory-Management Unit (MMU)	44
Dynamic relocation using a relocation register	45
Dynamic Linking	45
Swapping	45
Context Switch Time including Swapping	45
Swapping on Mobile Systems	46
Contiguous Memory Allocation	46
Dynamic Storage-Allocation Problem	47
Fragmentation	47
Segmentation	47
Segmentation Architecture	47
Paging	48
Address Translation Scheme	48
Implementation of Page Table	49
Memory Protection	49
Shared Pages	50
Structure of the Page Table	50
Hierarchical Page Tables	50
Hashed Page Tables	51
Inverted Page Table	51
Virtual Memory	51
Background	52
Virtual-address Space	52
Demand Paging	52
Page Fault	53
Aspects of Deman Paging	53
Performance of Demand Paging	54
Demand Paging Optimizations	54
Copy-on-Write	55
What happens if there is no free frame?	55
Page Replacement	55
Basic Page Replacement	55
Page and Frame Replacement Algorithms	55
First-In-First-Out (FIFO) Algorithm	55
Optimal Algorithm	56
Least Recently Used (LRU) Algorithm	56
LRU Approximation Algorithms	56
Enhanced Second-Chance Algorithm	57
Counting Algorithms	57
Page-Buffering Algorithms	57

Non-Uniform Memory Access	57
Thrashing	58
Demand Paging and Thrashing	58
Working-Set Model	58
Memory-Mapped Files	58
Allocating Kernel Memory	59
Buddy System	59
Slab Allocator	59
Other Issues	59
Page Size	59
TLB Reach	60
Program Structure	60
I/O interlock	60
Mass-Storage System	60
Overview of Mass Storage Structure	60
Magnetic Disks	61
Solid-State Disks	61
SSHD	61
Magnetic Tape	61
Disk Structure	61
Disk Attachment	62
Storage Array	62
Storage Area Network	62
Network-Attached Storage	63
Disk Scheduling	63
FCFS	63
SSTF	63
SCAN	63
C-SCAN	63
C-Look	64
Selecting a Disk-Scheduling Algorithm	64
Disk Management	64
Swap-Space Management	64
RAID Structure	65
Other Features	65
Extensions	65
Stable-Storage Implementation	65
File-System Interface	66
File Concept	66
File Attributes	66
File Operations	66

Close	67
Open File Locking	67
File Structure	67
Access Methods	67
Disk and Directory Structure	68
Directory Structure	68
Disk Structure	68
Operations Performed on Directory	68
Organize the Directory (Logically) to Obtain	68
File Sharing	69
Remote File Systems	69
Failure Modes	69
Protection	69
Access Lists and Groups	69
File System Implementation	70
File-System Structure	70
File System Layers	70
File-System Implementation	70
In-Memory File System Structures	71
Directory Implementation	71
Allocation Methods	71
Contiguous	72
Extent-Based Systems	72
Linked	72
Indexed	72
Performance	72
Free-Space Management	72
Efficiency and Performance	73
Recovery	73
Log Structured File Systems	73
NFS	73
NFS Protocol	74

List of Tables

Placeholder for table of contents	0
-----------------------------------	---

List of Figures

Figure 1: Process in Memory	17
Figure 2: Communications Models	19
Figure 3: Multithreaded Server Architecture	23
Figure 4: Single and Multithreaded Processes	23
Figure 5: CPU Burst	27

Figure 6: FCFS	29
Figure 7: Preemptive SJF Gantt Chart	29
Figure 8: Priority Scheduling Example	30
Figure 9: Resource-Allocation Graph	41
Figure 10: Multistep Processing of a User Program	44
Figure 11: Paging Hardware	48
Figure 12: Paging Hardware with TLB	49
Figure 13: Two-Level Page-Table Scheme	51
Figure 14: Steps in Handling a Page Fault	53
Figure 15: Optimal	56
Figure 16: LRU	56
Figure 17: Working-Set Model	58
Figure 18: Storage Area Network	62
Figure 19: A Typical File-system Organization	68
Figure 20: In-Memory File System Structures	71

Course Outline

- Intro
 - What is an OS? Major roles/responsibilities. User/kernel interaction. Operating system structures
- Processes and Threads
 - Operations on processes, Process state, Program vs process, threads vs processes, IPC
- Scheduling
 - Concepts, criteria, algorithms, threads and scheduling
- Deadlock and synchronization
 - Prevention/avoidance/detection/recovery
- Memory management and virtual memory
 - The memory hierarchy, swapping, paging
 - Demand paging, copy-on-write, page replacement, mem-mapped files
- I/O subsystems
 - IO hardware, device models, drivers, interrupt handling, DMA
- Mass storage and filesystems
 - Disks, file systems, mounting, network file systems, disk IO scheduling
- Specialized OSes
 - Real time systems, multimedia systems, embedded systems
- Protection and Security

Introduction

What is an operating system?

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

What Operating Systems Do

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

Operating System Definition

- OS is a resource allocator
 - Manages all resources
 - Decides between conflicting requests for efficient and fair resource use
- OS is a control program
 - Controls execution of programs to prevent errors and improper use of the computer
- No universally accepted definition
- "Everything a vendor ships when you order an operating system" is good approximation (But varies wildly)
- "The one program running at all times on the computer" is the kernel. Everything else is either a system program (ships with the operating system) or an application program

Common Functions of Interrupts

- Interrupt transfers control of the interrupt service routine generally, through the interrupt vector, which contains the addresses of all the service routines
- Interrupt architecture must save the address of the interrupted instruction
- A trap or exception is a software-generated interrupt caused either by an error or a user request
- An operating system is interrupt driven

Interrupt Handling

- The operating system preserves the state of the CPU by storing registers and the program counter
- Determines which type of interrupt has occurred:
 - polling
 - vectored interrupt system
- Separate segments of code determine what action should be taken for each type of interrupt

Computer-System Architecture

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

Operating System Structure

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

Operating-System Operations

- Interrupt driven by hardware
- Software error or requests creates exception or trap
 - Division by zero, request for operating system service
- Other process problems include infinite loop, processes modifying each other or the operating system
- Dual-mode operation allows OS to protect itself and other system components
 - User mode and kernel mode

-
- Mode bit provided by hardware
 - Provides ability to distinguish when system is running user code or kernel code
 - Some instructions designated as privileged, only executable in kernel mode
 - System call changes mode to kernel, return from call resets it to user
 - Increasingly CPUs support multi-mode operations
 - i.e. virtual machine manager (VMM) mode for guest VMs

Transition from User to Kernel Mode

- Timer to prevent infinite loop / process hogging resources
 - Set interrupt after specific period
 - Operating system decrements counter
 - When counter zero generate an interrupt
 - Set up before scheduling process to regain control or terminate program that exceeds allotted time

Process Management

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

Memory Management

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

Storage Management

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

IO Subsystem

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

Protection and Security

- Protection – any mechanism for controlling access of processes or users to resources defined by the OS
- Security – defense of the system against internal and external attacks
 - Huge range, including denial-of-service, worms, viruses, identity theft, theft of service

Open-Source Operating Systems

- Operating system made available in source-code format rather than just binary closed-source
- Counter to the copy protection and Digital Rights Management (DRM) movement
- Started by Free Software Foundation (FSF), which has "copyleft" GNU Public License (GPL)
- Examples include GNU/Linux and BSD UNIX (including core of Mac OS X), and many more

Operating-System Structures

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), 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)
 - **IO operations** – A running program may require IO, which may involve a file or an IO 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
 - **Communications** – Processes may exchange information, on the same computer or between computers over a network
 - **Error detection** – OS needs to be constantly aware of possible errors
 - May occur in the CPU and memory hardware, in IO 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
- Another set of OS functions exists for ensuring the efficient operation of the system itself via resource sharing
 - **Resource allocation** – When multiple users or multiple jobs running concurrently, resources must be allocated to each of them
 - Many types of resources. Some (such as CPU cycles, main memory, and file storage) may have special allocation code, others (such as IO devices) may have general request and release code
 - **Accounting** – To keep track of which users use how much and what kinds of computer resources
 - **Protection and security** – The owners of information stored in a multiuser or networked computer system may want to control use of that information, concurrent processes should not interfere with each other
 - **Protection** involves ensuring that all access to system resources is controlled
 - **Security** of the system from outsiders requires user authentication, extends to defending external IO devices from invalid access attempts
 - If a system is to be protected and secure, precautions must be instituted throughout it. A chain is only as strong as its weakest link

User Operating System Interface

CLI

CLI or command interpreter 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

GUI

- User-friendly desktop metaphor interface
 - Usually mouse, keyboard, and monitor
 - Icons represent files, programs, actions, etc
 - Various mouse buttons over objects in the interface cause various actions (provide information, options, execute function, open directory)
- Many systems now include both CLI and GUI interfaces

Touchscreen

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

System Calls

System Call Parameter Passing

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

Types of System Calls

- Process control
 - end, abort
 - load, execute
 - create process, terminate process
 - 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
- 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
 - 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
 - Protection
 - control access to resources
 - get and set permissions
 - allow and deny user access

System Programs

- System programs provide a convenient environment for program development and execution. They can be divided into:
 - File manipulation
 - Status information sometimes stored in a File modification
 - Programming language support
 - Program loading and execution
 - Communications
 - Background services
 - Application programs
- Most users' view of the operation system is defined by system programs, not the actual system calls
- Provide a convenient environment for program development and execution
 - Some of them are simply user interfaces to system calls; others are considerably more complex
- **File management** – Create, delete, copy, rename, print, dump, list, and generally manipulate files and directories
- **Status information**
 - Some ask the system for info – date, time, amount of available memory, disk space, number of users
 - Others provide detailed performance, logging, and debugging information
 - Typically, these programs format and print the output of the terminal or other output devices
 - Some systems implement a registry – used to store and retrieve configuration information
- **File modification**
 - Text editors to create and modify files
 - Special commands to search contents of files or perform transformations of the text
- **Programming-language support** – Compilers, assemblers, debuggers and interpreters sometimes provided
- **Program loading and execution** – Absolute loaders, relocatable loaders, linkage editors, and overlay-loaders, debugging systems for higher-level and machine language
- **Communications** – Provide the mechanism for creating virtual connections among processes, users, and computer systems
 - Allow users to send messages to one another's screens, browse web pages, send electronic-mail messages, log in remotely, transfer files from one machine to another
- **Background Services**

-
- Launch at boot time
 - Some for system startup, then terminate
 - Some from system boot to shutdown
 - Provide facilities like disk checking, process scheduling, error logging, printing
 - Run in user context not kernel context
 - Known as services, subsystems, daemons
 - **Application programs**
 - Don't pertain to system
 - Run by users
 - Not typically considered part of OS
 - Launched by command line, mouse click, finger poke

Operating System Design and Implementation

- Design and Implementation of best OS not "solvable", but some approaches have proven successful
- User goals and System goals
 - **User goals** – operating system should be convenient to use, easy to learn, reliable, safe, and fast
 - **System goals** – operating system should be easy to design, implement, and maintain, as well as flexible, reliable, error-free, and efficient
- Important principle to separate
 - > Policy: What will be done?
 - > Mechanism: How to do it?
- Mechanisms determine how to do something, policies decide what will be done
 - The separation of policy from mechanism is a very important principle, it allows maximum flexibility if policy decisions are to be changed later
- Specifying and designing OS is highly creative task of software engineering

Implementation

- Much variation
 - Early OSes in assembly language
 - Then system programming languages like Algol, PL/1
 - Now C, C++
- Actually usually a mix of languages
 - Lowest levels in assembly
 - Main body in C
 - System programs in C, C++, scripting languages like PERL, Python, shell scripts
- More high-level language easier to port to other hardware (but slower)
- Emulation can allow an OS to run on non-native hardware

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

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

Microkernel System Structure

- Moves as much from the kernel into user space
- Communication takes place between user modules using message passing
- Benefits:
 - Easier to extend a microkernel
 - Easier to port the operating system to new architectures
 - More reliable (less code is running in kernel mode)
 - More secure
- Detriments:
 - Performance overhead of user space to kernel space communication

Modules

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

System Boot

- When power initialized on system, execution starts at a fixed memory location
 - Firmware ROM used to hold initial boot code
- Operating system must be made available to hardware so hardware can start it
 - Small piece of code – bootstrap loader, stored in ROM or EEPROM locates the kernel, loads it into memory, and starts it
 - Sometimes two-step process where boot block at fixed location loaded by ROM code, which loads bootstrap loaded from disk
- Common bootstrap loader, GRUB, allows selection of kernel from multiple disks, versions, kernel options
- Kernel loads and system is then running

Processes

Process Concept

- An operating system executes a variety of programs
 - Batch system – [jobs](#)

Time-sharing systems – [user programs](#) or [tasks](#)

- Textbook uses the terms *job* and *process* almost interchangeably
- **Process** – a program in execution; process execution must progress in sequential fashion
- Multiple parts
 - The program code, also call [text section](#)
 - Current activity including [program counter](#), processor registers
 - **Stack** containing temporary data
 - Function parameters, return addresses, local variables
 - **Data section** containing global variables
 - **Heap** containing memory dynamically allocated during run time
- Program is *passive* entity stored on disk ([executable file](#)), process is *active*

- Program becomes process when executable file 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

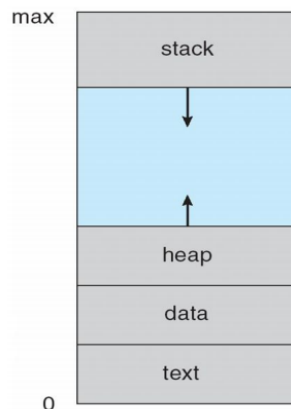


Figure 1: Process in Memory

Process State

- As a process executes, it changes **state**
 - **New:** The process is being created
 - **Running:** Instructions are being executed
 - **Waiting:** The process is waiting for some event to occur
 - **Ready:** The process is waiting to be assigned to a processor
 - **Terminated:** The process has finished execution

! [Diagram of Process State](sem2-2017/comp3301/state.png)₇₅

Process Control Block (PCB)

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

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

Threads

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

Process Scheduling

- Maximize CPU use, quickly switch processes onto CPU for time sharing
- **Process scheduler** selects among available processes for next execution on CPU
- Maintains **scheduling queues** of processes
 - **Job queue** – set of all processes in the system
 - **Ready queue** – set of all processes residing in main memory, ready and waiting to execute
 - **Device queues** – set of processes waiting for an I/O device

-
- Processes migrate among the various queues

Schedulers

- **Long-term scheduler** (or [job scheduler](#)) – selects which processes should be brought into the ready queue
- **Short-term scheduler** (or [CPU scheduler](#)) – selects which process should be executed next and allocates CPU
 - Sometimes the only scheduler in a system
- Short-term scheduler is invoked very frequently (milliseconds) → (must be fast)
- Long-term scheduler is invoked very infrequently (seconds, minutes) → (may be slow)
- The long-term scheduler controls the [degree of multiprogramming](#)
- Processes can be described as either:
 - **I/O-bound process** – spends more time doing I/O than computations, many short CPU bursts
 - **CPU-bound process** – spends more time doing computations; few very long CPU bursts
- Long-term scheduler strives for good *process mix*

Addition of Medium Term Scheduling

- **Medium-term scheduler** can be added if degree of multiple programming needs to decrease
 - Remove process from memory, store on disk, bring back in from disk to continue execution → [swapping](#)

Multitasking in Mobile Systems

- Some systems / early systems allow only one process to run, other suspended
- Due to screen real estate, user interface limits iOS provides for a
 - Single **foreground** process – controlled via user interface
 - Multiple **background** processes – in memory, running, but not on the display, and with limits
 - Limits include single, short task, receiving notification of events, specific long-running tasks like audio playback
- Android runs foreground and background, with fewer limits
 - Background process uses a [service](#) to perform tasks
 - Service can keep running even if background process is suspended
 - Service has no user interface, small memory use

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 process via a [context switch](#)
- **Context** of a process represented in the PCB
- Context-switch time is overhead; the system does no useful work while switching
 - The more complex the OS and the PCB → 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, termination, and so on as detailed

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

Process Termination

- Process executes last statement and asks the operating system to delete it (`exit()`)
 - Output data from child to parent (via `wait()`)
 - Process' resources are deallocated by operating system
- Parent may terminate execution of children processes (`abort()`)
 - Child has exceeded allocated resources
 - Task assigned to child is no longer required
 - If parent is exiting
 - Some operating systems do not allow child to continue if its parent terminates
 - All children terminated – [cascading termination](#)
- Wait for termination, returning the pid:


```
1 | pid t_pid; int status;
2 | pid = wait(&status);
```
- If no parent waiting, then terminated process is a **zombie**
- If parent terminated, processes are **orphans**

Interprocess Communication

- Processes within a system may be *independent* or *cooperating*
- Cooperating process can affect or be affected by other processes, including sharing data
- Reasons for cooperating processes:
 - Information sharing
 - Computation speedup
 - Modularity
 - Convenience
- Cooperating processes need [interprocess communication \(IPC\)](#)
- Two models of IPC
 - **Shared memory**
 - **Message passing**

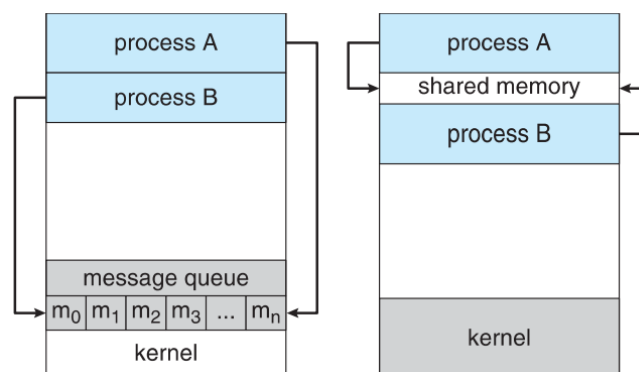


Figure 2: Communications Models

Cooperating Processes

-
- *Independent* process cannot affect or be affected by the execution of another process
 - *Cooperating* process can affect or be affected by the execution of other process
 - Advantages of process cooperation
 - Information sharing
 - Computation speed-up
 - Modularity
 - Convenience

Producer-Consumer Problem

- Paradigm for cooperating processes, *producer* process produces information that is consumed by a *consumer* process
 - **unbounded-buffer** places no practical limit of the size of the buffer
 - **bounded-buffer** assumes that there is a fixed buffer size

Interprocess Communication - Message Passing

- Mechanism for processes to communicate and to synchronize their actions
- Message system – processes communicate with each other without resorting to shared variables
- IPC facility provides two operations:
 - **send(message)** – message size fixed or variable
 - **receive(message)**
- If *P* and *Q* wish to communicate, they need:
 - establish a **communication link** between them
 - exchange message via send/receive
- Implementation of communication link
 - physical (e.g. shared memory, hardware bus)
 - logical (e.g. direct or indirect, synchronous or asynchronous, automatic or explicit buffering)

Direct Communication

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

Indirect Communication

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

-
- Primitives are defined as:
 - `send(A, message)` – send a message to mailbox A
 - `receive(A, message)` – receive a message from mailbox A

Synchronization

- Message passing may be either blocking or non-blocking
- **Blocking** is considered **synchronous**
 - Blocking send has the sender block until the message is received
 - Blocking receive has the receiver block until a message is available
- **Non-blocking** is considered **asynchronous**
 - Non-blocking send has the sender send the message and continue
 - Non-blocking receive has the receiver receive a valid message or null
- Different combinations possible
 - If both send and receive are blocking, we have a **rendezvous**
- Producer-consumer becomes trivial

Buffer

- Queue of messages attached to the link; implemented in one of three ways
 - 1) Zero capacity – 0 messages
Sender must wait for receiver (rendezvous)
 - 1) Bounded capacity – finite length of n messages
Sender must wait if link full
 - 1) Unbounded capacity – infinite length
Sender never waits

Examples of IPC Systems

POSIX

- POSIX Shared Memory
 - Process first creates shared memory segment

```
shm_fd = shm_open(name, O_CREAT | O_RDWR, 0666);
```
 - Also used to open an existing segment to share it
 - Set the size of the object

```
ftruncate(shm_fd, 4096);
```
 - Now the process could write to the shared memory

```
sprintf(shared_memory, "Writing to shared memory");
```

Communication in Client-Server Systems

Sockets

- A **socket** is defined as an endpoint for communication
- Concatenation of IP address and **port** – a number included at start of message packet to differentiate network services on a host
- The socket **161.25.19.8:1625** refers to port **1625** on host **161.25.19.8**
- Communication consists between a pair of sockets
- All ports below 1024 are *well known*, used for standard services
- Special IP address **127.0.0.1** (**loopback**) to refer to system on which process is running

Remote Procedure Calls

- Remote procedure call (RPC) abstracts procedure calls between processes on networked systems
 - Again uses ports for service differentiation
- **Stubs** – client-side proxy for the actual procedure on the server
- The client-side stub locates the server and **marshalls** the parameters

- The server-side stub receives this message, unpacks the marshalled parameters, and performs the procedure on the server
- On Windows, stub code compile from specification written in **Microsoft Interface Definition Language (MIDL)**
- Data representation handled via **External Data Representation (XDL)** format to account for different architectures
 - [Big-endian](#) and [little-endian](#)
- Remote communication has more failure scenarios than local
 - Messages can be delivered *exactly once* rather than *at most once*
- OS typically provides a rendezvous (or [matchmaker](#)) service to connect client and server

Pipes

- Acts as a conduit allowing two processes to communicate
- Issues
 - Is communication unidirectional or bidirectional?
 - In the case of two-way communication, is it half or full-duplex?
 - Must there exist a relationship (i.e. parent-child) between the communicating processes?
 - Can the pipes be used over a network?

Ordinary Pipes

- Ordinary Pipes allow communication in standard producer-consumer style
 - Producer writes to one end (the [write-end](#) of the pipe)
 - Consumer reads from the other end (the [read-end](#) of the pipe)
 - Ordinary pipes are therefore unidirectional
 - Require parent-child relationship between communicating processes
- Windows calls these anonymous pipes*

Named Pipes

- Named Pipes are more powerful than ordinary pipes
- Communication is bidirectional
- No parent-child relationship is necessary between the communicating processes
- Several processes can use the named pipe for communication
- Provided on both UNIX and Windows systems

Threads

Overview

Motivation

- Most modern applications are multithreaded
- Threads run within application
- Multiple tasks with the application can be implemented by separate threads
 - Update display
 - Fetch data
 - Spell checking
 - Answer a network request
- Process creation is heavy-weight while thread creation is light-weight
- Can simplify code, increase efficiency
- Kernels are generally multithreaded

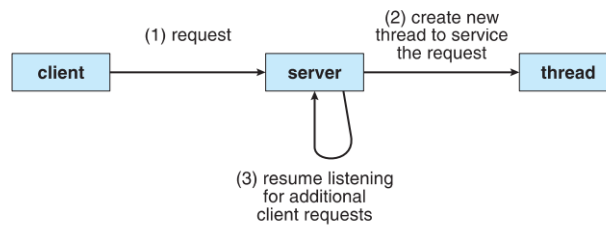


Figure 3: Multithreaded Server Architecture

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 multiprocessor architectures

Multicore Programming

- **Multicore** or **multiprocessor** systems putting 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
- 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
- As number of threads grows, so does architectural support for threading
 - CPUs have cores as well as *hardware threads*

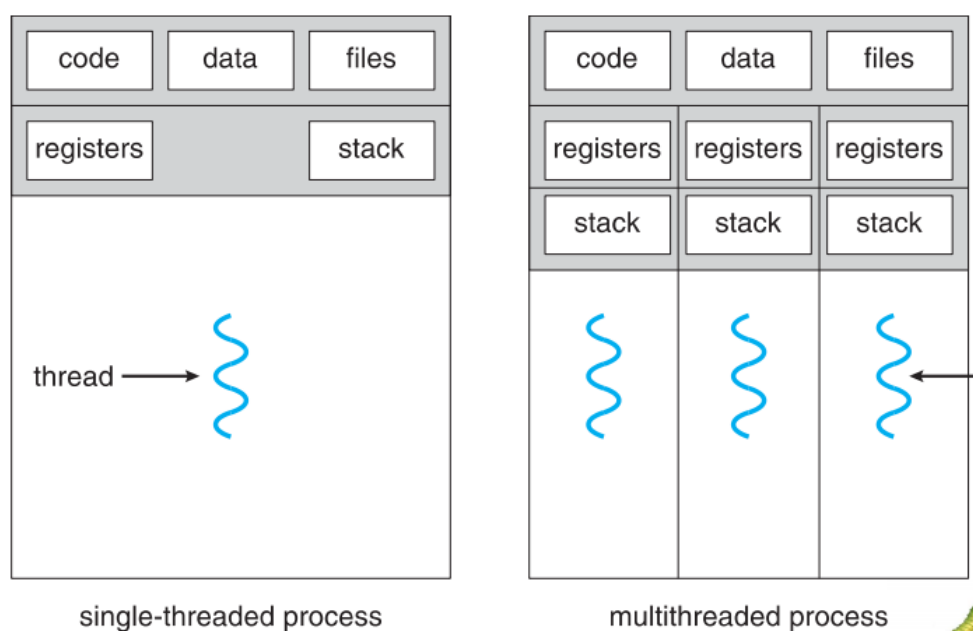


Figure 4: Single and Multithreaded Processes

Amdahls Law

- Identifies performance gains from adding additional cores to an application that has both serial and parallel components
- S is a serial portion
- N processing cores

$$speedup \leq \frac{1}{S + \frac{(1-S)}{N}}$$

- i.e. if application is 75% parallel / 25% serial, moving from 1 to 2 cores results in speedup of 1.6 times
 - As N approaches infinity, speedup approaches $1/S$
- > Serial portion of an application has disproportionate effect on performance gained by adding additional cores
- But does the law take into account contemporary multicore systems?

User Threads and Kernel Threads

- **User threads** – management done by user-level threads library
- Three primary thread libraries
 - 1) POSIX pthreads
 - 2) Windows threads
 - 3) Java threads
- Kernel threads – supported by the Kernel
- Examples – virtually all general purpose operating systems, including:
 - Windows, Solaris, Linux, Tru64 UNIX, Mac OS X

Multithreading Models

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

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

Many-to-Many

- Allows many user level threads to be mapped to many kernel threads
- Allows the operating system to create a sufficient number of kernel threads

Thread Libraries

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

Pthreads

- 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

Java Threads

- 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

Implicit Threading

- Growing in popularity as numbers of threads increase, program correctness more difficult with explicit threads
- Creation and management of threads done by compilers and run-time libraries rather than programmers
- Three methods explored
 - Thread Pools
 - OpenMP
 - Grand Central Dispatch
- Other methods include Microsoft Threading Building Blocks (TBB), `java.util.concurrent` package

Thread Pools

- Create a number of threads in a pool where they await work
 - Advantages:
 - Usually slightly faster to service a request with an existing thread than create a new thread
 - Allows the number of threads in the application(s) to be bound to the size of the pool
 - Separating task to be performed from mechanics of creating task allows different strategies for running task
- > i.e. Tasks could be scheduled to run periodically

OpenMP

- Set of compiler directives and an API for C, C++, FORTRAN
- Provides support for parallel programming in shared-memory environments
- Identifies [parallel regions](#) – blocks of code that can run in parallel

Create as many threads as there are cores _____

```
...  
#pragma omp parallel  
...
```

Run for loop in parallel _____

```
...  
#pragma omp parallel for  
for (i=0; i<N; i++) {  
    c[i] = a[i] + b[i];  
}  
...
```

Grand Central Dispatch

- Apple technology for Mac OS X and iOS operating systems
- Extensions to C, C++ languages, API, and run-time library
- Allows identification of parallel sections
- Manages most of the details of threading
- Block is in `^{} – ^{ printf("I am a block"); }`
- Blocks placed in dispatch queue
 - Assigned to available thread in thread pool when removed from queue
- Two types of dispatch queues:
 - **Serial** – blocks removed in FIFO order, queue is per process, called [main queue](#)

- Programmers can create additional serial queues within program
- **Concurrent** – removed in FIFO order but several may be removed at a time
 - Three system wide queues with priorities low, default, high

```
1 | dispatch_queue_t queue = dispatch_get_global_queue(DISPATCH_QUEUE_PRIORITY_DEFAULT, 0);
2 | dispatch_async(queue, ^{ printf("I am a block"); });
```

Threading Issues

- 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()`

- 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

Signal Handling

- **Signals** are used in UNIX systems to notify a process that a particular event has occurred
- A **signal handler** is used to process signals
 - 1) Signal is generated by particular event
 - 2) Signal is delivered to a process
 - 3) Signal is handled by one of two signal handlers:
 - 1) default
 - 2) user-defined
- Every signal has **default handler** that kernel runs when handling signal
 - **User-defined signal handler** can override default
 - For single-threaded, signal delivered to process
- Where should a signal be delivered for multi-threaded?
 - Deliver the signal to the thread to which the signal applies
 - Deliver the signal to every thread in the process
 - Deliver the signal to certain threads in the process
 - Assign a specific thread to receive all signals for the process

Thread-Local Storage

- **Thread-local storage (TLS)** allows each thread to have its own copy of data
- Useful when you do not have control over the thread creation process (i.e. when using a thread pool)
- Different from local variables
 - Local variables visible only during single function invocation
 - TLS visible across function invocations
- Similar to `static` data
 - TLS is unique to each thread

Linux Threads

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

- Flags control behavior

Flag

Meaning

`CLONE_FS`

File-system information is shared

`CLONE_VM`

The same memory space is shared

`CLONE_SIGHAND`

Signal handlers are shared

`CLONE_FILES`

The set of open files is shared

- `struct task_struct` points to process data structures (shared or unique)

CPU Scheduling

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

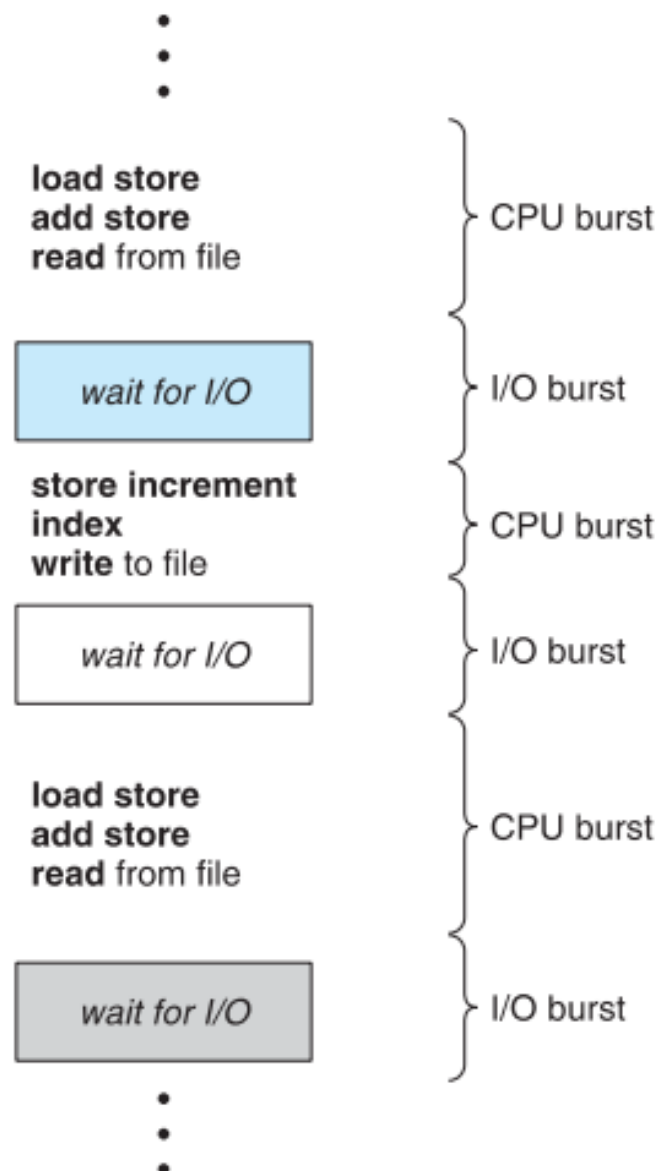


Figure 5: CPU Burst

CPU Scheduler

- **Short-term scheduler** selects from among the processes in ready queue, and allocates the CPU 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
- Scheduling under 1 and 4 is **nonpreemptive**
- All other scheduling is **preemptive**
 - Consider access to shared data
 - Consider preemption while in kernel mode
 - Consider interrupts occurring during crucial OS activities

Dispatcher

- Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:
 - Switching context
 - Switching to user mode
 - Jumping to the proper location in the user program to restart that program
- **Dispatch latency** – time it takes for the dispatcher to stop one process and start another running

Scheduling Criteria

- **CPU utilization** – keep the CPU as busy as possible
- **Throughput** – number 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, (for time-sharing environment)

Scheduling Algorithms

Criteria

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

> May be interested in AVERAGE or WORST CASE figures

First-Come, First-Served (FCFS) Scheduling

Process	Burst Time
P ₁	24
P ₂	3
P ₃	3

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

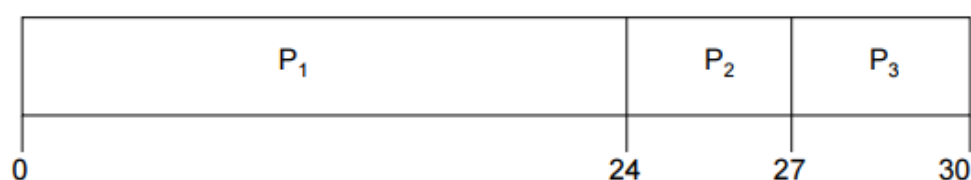


Figure 6: FCFS

- Average waiting time is 17
- Suppose that the processes arrive in the order: P_2, P_3, P_1
- Average waiting time is 3
 - Much better than previous case
 - **Convoy effect** – short process behind long process
 - Consider one CPU-bound and many I/O-bound processes

Shortest-Job-First (SJF) Scheduling

- Associate with each process the length of its next CPU burst
 - Use these lengths to schedule the process with the shortest time
- SJF is optimal – gives minimum average waiting time for a given set of processes
 - The difficulty is knowing the length of the next CPU request
 - Could ask the user

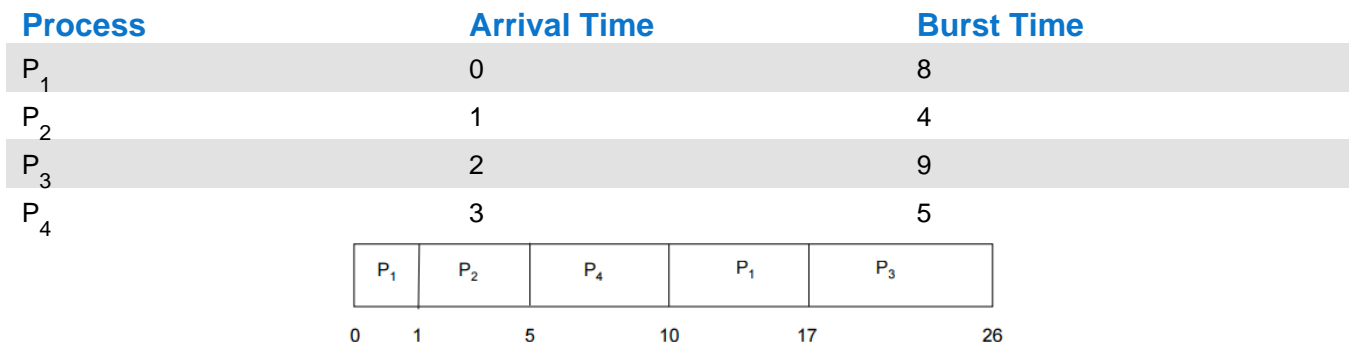
Determining Length of Next CPU Burst

- Can only estimate the length – should be similar to the previous one
- Then pick process with shortest predicted next CPU burst
- Can be done by using the length of previous CPU bursts, using exponential averaging
 1. t_n = actual length of n^{th} CPU burst
 1. $\text{predicted value for the next CPU burst}$
 1. $\alpha, 0 \leq \alpha \leq 1$
 1. Define:

$$\mathcal{t}_{n+1} = \alpha t_n + (1 - \alpha) \mathcal{t}_n$$
 - Commonly, α set to 1/2
 - Preemptive version called **shortest-remaining-time-first**

Example of Shortest-remaining-time-first

- Now we add the concepts of varying arrival times and preemption to the analysis



- Average waiting time = 6.5

Priority Scheduling

- A priority number (integer) is associated with each process
- The CPU is allocated to the process with the highest priority (smallest integer = highest priority)
 - Preemptive
 - Nonpreemptive
- SJF is priority scheduling where priority is the inverse of predicted next CPU burst time
- Problem = **Starvation** – low priority processes may never execute
- Solution = **Ageing** – as time progresses increase the priority of the process

Example of Priority Scheduling

Process	Burst Time	Priority
P ₁	10	3
P ₂	1	1
P ₃	2	4
P ₄	1	5
P ₅	5	2

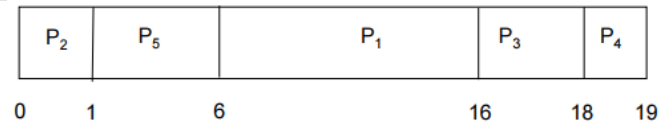


Figure 8: Priority Scheduling Example

- Average waiting time = 8.2

Round Robin

- Each process gets a small unity of CPU time (**time quantum** q), usually 10-100 milliseconds. After this time has elapsed, the process is preempted and added to the end of the ready queue
- If there are n processes in the ready queue and the time quantum is q , then each process gets $1/n$ of the CPU time in chunks of at most q time units at once. No process waits more than $(1-n)q$ time units.
- Timer interrupts every quantum to schedule next process
- Performance
 - q large \rightarrow FIFO
 - q small $\rightarrow q$ must be large with respect to context switch, otherwise overhead is too high

Multilevel Queue

- Ready queue is partitioned into separate queues, e.g.
 - **foreground** (interactive)
 - **background** (batch)
- Process permanently in a given queue
- Each queue has its own scheduling algorithm:
 - foreground – RR
 - background – FCFS
- Scheduling must be done between the queues:
 - Fixed priority scheduling (i.e. serve all from foreground then from background). Possibility of starvation
 - Time slice – each queue gets a certain amount of CPU time which it can schedule amongst its processes; i.e. 80% to foreground in RR
 - 20% to background in FCFS

Multilevel Feedback Queue

- A process can move between the various queue; aging can be implemented this way
- Multilevel-feedback-queue scheduler defined by the following parameters:
 - number of queues
 - scheduling algorithms for each queue (may be different)
 - 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
- Most general algorithm, but also most complicated

Thread Scheduling

- Distinction between user-level and kernel-level threads
- When threads supported, threads scheduled, not processes

-
- Many-to-one and many-to-many models, thread library schedules user-level threads to run on LWP (lightweight processes)
 - Known as **process-contention scope (PCS)** since scheduling competition is within the process
 - Typically done via priority set by programmer
 - Kernel thread scheduled onto available CPU is **system-contention scope (SCS)** – competition among all threads in system

Multiple-Processor Scheduling

- CPU scheduling more complex when multiple CPUs are available
- **Homogeneous processors** within a multiprocessor
- **Asymmetric multiprocessing (SMP)** – each processor is self-scheduling, all processes in common ready queue, or each has its own private queue of ready processes
 - Currently, most common
- **Processor affinity** – process has affinity for processor on which it is currently running
 - **soft affinity**
 - **hard affinity**
 - Variations including processor sets

Load Balancing

- 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

Multicore Processors

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

Virtualization and Scheduling

- Virtualization software schedules multiple guests onto CPU(s)
- Each guest doing its own scheduling
 - Not knowing it doesn't own the CPUs
 - Can result in poor response time
 - Can effect time-of-day clocks in guests
- Can undo good scheduling algorithm efforts of guests

Real-Time CPU Scheduling

- Can present obvious challenges
- **Soft real-time systems** – no guarantee as to when critical real-time process will be scheduled
- **Hard real-time systems** – task must be serviced by its deadline
- 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 schedule to take current process off CPU and switch to another
- Conflict phase of dispatch latency:
 - 1) Preemption of any process running in kernel mode
 - 2) Release by low-priority process of resources

Priority-based Scheduling

- 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 $1/p$

Rate Monotonic Scheduling

- 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

Earliest Deadline First Scheduling (EDF)

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

Proportional Share Scheduling

- 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 Systems Examples

Linux Scheduling

- **Completely Fair Scheduler** (CFS)
- **Scheduling classes**
 - Each has specific priority
 - Scheduler picks highest priority task in highest scheduling class
 - Rather than quantum based on fixed time allotments, based on proportion of CPU time
 - 2 scheduling classes included, others can be added
 - 1) default
 - 2) real-time
- Quantum calculated based on **nice value** from -20 to +19
 - Lower value is higher priority
 - Calculates **target latency** – interval of time during which task should run at least once
 - Target latency can increase if say number of active tasks increases
- CFS scheduler maintains per task **virtual run time** in variable **vruntime**
 - Associated with decay factor based on priority of task – lower priority is higher decay rate
 - Normal default priority yields virtual run time = actual run time
- To decide next task to run, scheduler picks task with lowest virtual run time
- Real-time scheduling according to POSIX.1b
 - Real-time tasks have static priorities
- Real-time plus normal map into global priority scheme
- Nice value of -20 maps to global priority 100
- Nice value of +19 maps to priority 139

Windows Scheduling

- Windows uses priority-based preemptive scheduling

- Highest-priority thread runs next
- **Dispatcher** is scheduler
- Thread runs until:
 - 1) blocks
 - 2) uses time slice
 - 3) preempted by higher-priority thread
- Real-time threads can preempt non-real-time
- 32-level priority scheme
- **Variable class** is 1-15, **real-time class** is 16-31
- Priority 0 is memory-management thread
- Queue for each priority
- If no run-able thread, runs **idle thread**

Solaris Scheduling

- Priority-based scheduling
- Six classes available
 - Time sharing (default) (TS)
 - Interactive (IA)
 - Real time (RT)
 - System (SYS)
 - Fair Share (FSS)
 - Fixed priority (FP)
- Given thread can be in one class at a time
- Each class has its own scheduling algorithm
- Time sharing is multi-level feedback queue
 - Loadable table configurable by sysadmin

Algorithm Evaluation

- How to select CPU-scheduling algorithm for an OS?
- Determine criteria, then evaluate algorithms
- **Deterministic modelling**
 - Type of [analytic evaluation](#)
 - Takes a particular predetermined workload and defines the performance of each algorithm for that workload
- Consider 5 processes arriving at time 0

Process	Burst Time
P ₁	10
P ₂	29
P ₃	3
P ₄	7
P ₅	12

Deterministic Evaluation

- For each algorithm, calculate minimum average waiting time
- Simple and fast, but requires exact numbers for input, applies only to those inputs
 - FCS is 28ms
 - Non-preemptive SFJ is 13ms
 - RR is 23ms

Queueing Models

-
- Describes the arrival of processes, and CPU and IO bursts probabilistically
 - Commonly exponential, and described by mean
 - Computes average throughput, utilization, waiting time, etc
 - Computer system described as network of servers, each with queue of waiting processes
 - Knowing arrival rates and service rates
 - Computes utilization, average queue length, average wait time, etc

Little's Law

- n = average queue length
- W = average waiting time in queue
- λ = average arrival rate into queue
- Little's law - in steady state, processes leaving queue must equal processes arriving, thus $n = \lambda \times W$
 - Valid for any scheduling algorithm and arrival distribution
- For example, if on average 7 processes arrive per second, and normally 14 processes in queue, then average wait time per process = 2 seconds

Simulations

- Queueing models limited
- **Simulations** more accurate
 - Programmed model of computer system
 - Clock is a variable
 - Gather statistics indicating algorithm performance
 - Data to drive simulation gathered via
 - Random number generator according to probabilities
 - Distributions defined mathematically or empirically
 - Trace tapes record sequences of real events in real systems

Implementation

- Even simulations have limited accuracy
- Just implement new scheduler and test in real systems
 - High cost, high risk
 - Environments vary
- Most flexible schedulers can be modified per-site or per-system
- Or APIs to modify priorities
- But again environments vary

Process Synchronization

- To introduce the critical-section problem, whose solutions can be used to ensure the consistency of shared data
- To present both software and hardware solutions of the critical-section problem
- To examine several classical process-synchronization problems
- To explore several tools that are used to solve process synchronization problems

Background

- Processes can execute concurrently
 - May be interrupted at any time, partially completing execution
- Concurrent access to shared data may result in data consistency
- Maintaining data consistency requires mechanisms to ensure the orderly execution of cooperating processes

The Critical-Section Problem

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

Critical Section

General structure for a process:

```

1 | do {
2 |     entry section
3 |     critical section
4 |     exit section
5 |     remainder section
6 | } while (true);

```

Solution to Critical-Section Problem

- 1) **Mutual Exclusion** – If process p_i is executing in its critical section, then no other processes can be executing in their critical sections
 - 2) **Progress** – If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the processes that will enter the critical section next cannot be postponed indefinitely
 - 3) **Bounded Waiting** – A bound must exist on the number of times that other processes are allowed to enter their critical sections after a process has made a request to enter its critical section and before that request is granted
 - Assume that each process executes at a nonzero speed
 - No assumption concerning **relative speed** of the n processes
- Two approaches depending on if kernel is preemptive or non-preemptive
 - **Preemptive** – allows preemption of process when running in kernel mode
 - **Non-preemptive** – runs until exits kernel mode, blocks, or voluntarily yields CPU
 - Essentially free of race conditions in kernel mode

Petersons Solution

- Good algorithmic description of solving the problem
- Two process solution
- Assume that the **load** and **store** instructions are atomic; that is, cannot be interrupted
- The two processes share two variables:
 - **int turn;**
 - **bool flag[2];**
- The variable **turn** indicates whose turn it is to enter the critical section
- The **flag** array is used to indicate if a process is ready to enter the critical section. **flag[i] = true** implies that process p_i is ready

Algorithm for Process P_i

```

1 | do {
2 |     flag[i] = true;
3 |     turn = j;
4 |     while (flag[j] && turn == j);
5 |     critical section
6 |     flag[i] = false;
7 |     remainder section
8 | } while (true);

```

- Provable that

-
- 1) Mutual exclusion is preserved
 - 2) Progress requirement is satisfied
 - 3) Bounded-waiting requirement is met

Synchronization Hardware

- Many systems provide hardware support for critical section code
 - All solutions below based on idea of **locking**
 - Protecting critical regions via locks
 - Uniprocessors – could disable interrupts
 - Currently running code would execute without preemption
 - Generally too inefficient on multiprocessor systems
 - Operating systems using this not broadly scalable
 - Modern machines provide special atomic hardware instructions
- > Atomic = non-interruptible
- Either test memory word and set value
 - Or swap contents of two memory words

Mutex Locks

- Previous solutions are complicated and generally inaccessible to application programmers
- OS designers build software tools to solve critical section problems
- Simplest is mutex lock
- Protect critical regions with it by first `acquire()` a lock then `release()` it
 - Boolean variable indicating if lock is available or not
- Calls to `acquire()` and `release()` must be atomic
 - Usually implemented via hardware atomic instructions
- But this solution requires **busy waiting**
 - This lock therefore called a **spinlock**
 - Can be OK for short waits on a multi-processor system

Semaphore

- Synchronization tool that does not require busy waiting
- Semaphore S – integer variable
- Two standard operations modify S : `wait()` and `signal()`
 - Originally called `P()` and `V()`. The inventor of semaphores was Edsger Dijkstra who was very Dutch. **P** and **V** are the first letters of two Dutch words *proberen* (to test) and *verhogen* (to increment).
- Less complicated
- Can only be accessed via two indivisible (atomic) operations

Semaphore Usage

- **Counting semaphore** – integer value can range over an unrestricted domain
- **Binary semaphore** – integer value can range only between 0 and 1
 - Then a **mutex lock**
- Can implement a counting semaphore S as a binary semaphore
- Can solve various synchronization problems
- Consider P_1 and P_2 that require S_1 to happen before S_2

Semaphore Implementation

- Must guarantee that no two processes can execute `wait()` and `signal()` on the same semaphore at the same time

-
- Thus, implementation becomes the critical section problem where the wait and signal code are placed in the 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

Semaphore Implementation with no Busy waiting

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

Deadlock and Starvation

- **Deadlock** – two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes
- **Starvation** – [indefinite blocking](#)
 - A process may never be removed from the semaphore queue in which it is suspended
- **Priority Inversion** – Scheduling problem when lower-priority process holds a lock needed by higher-priority process
 - Solved via [priority-inheritance protocol](#)

Classic Problems of Synchronization

- Classical problems used to test newly-proposed synchronization schemes

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
- The structure of the producer process

```
1  do {
2      // produce an item in next_produced
3
4      wait(empty);
5      wait(mutex);
6
7      // add next produced to the buffer
8
9      signal(mutex);
10     signal(full);
11 } while (true);
```
- The structure of the consumer process

```

1 | do {
2 |     wait(full);
3 |     wait(mutex);
4 |
5 |     // remove an item from buffer to next_consumed
6 |
7 |     signal(mutex);
8 |     signal(empty);
9 |
10 |    // consume the item in next consumed
11 | } while (true);

```

Readers-Writers Problem

- A data set is shared among a number of concurrent processes
 - Readers – only read the data set; they do *not* perform any updates
 - Writers – can both read and write
- Problem – allow multiple readers to read at the same time
 - Only one single writer can access the shared data at the same time
- Several variations of how readers and writers are treated – all involve priorities
- Shared Data
 - Data set

Readers-Writers Problem Variations

- First variation – no reader kept waiting unless writer has permission to use shared object
- Second variation – once writer is ready, it performs write asap
- Both may have starvation leading to even more variations
- Problem is solved on some systems by kernel providing reader-writer locks

Dining-Philosophers Problem

- Philosophers spend their lives thinking and eating
- Don't interact with their neighbours, 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
 - Shared data
 - Bowl of rice (data set)
 - Semaphore `chopstick[5]` initialized to 1

Problems with Semaphores

- Incorrect use of semaphore operations:
 - `signal (mutex) ... wait (mutex)`
 - `wait (mutex) ... wait (mutex)`
 - Omitting of `wait(mutex)` or `signal(mutex)` (or both)
- Deadlock and starvation

Monitors

- 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 t a time
- But not powerful enough to model some synchronization schemes

```

1 | monitor monitor-name {
2 |     // shared variable declarations
3 |     procedure P1 (...) {
4 |         ...
5 |     }
6 |
7 |     procedure P2 (...) {
8 |         ...
9 |     }
10 |
11 |     Initialization code (...) {
12 |         ...
13 |     }
14 | }

```

Condition Variables

- **condition `x, y`;**
- Two operations on a condition variable:
 - **`x.wait()`** – a process that invokes the operation is suspended until **`x.signal()`**
 - **`x.signal()`** – resumes one of processes (if any) that invoked **`x.wait()`**
 - If no **`x.wait()`** on the variable, then it has no effect on the variable

Condition Variables Choices

- If process *P* invokes **`x.signal()`**, with *Q* in **`x.wait()`** state, what should happen next?
 - If *Q* is resumed, then *P* must wait
- Options include
 - **Signal and wait** – *P* waits until *Q* leaves monitor or waits for another condition
 - **Signal and continue** – *Q* waits until *P* leaves the monitor or waits for another condition
- Both have pros and cons – language implmenter can decide
- Monitors implemented in Concurrent Pascal compromise
 - *P* executing signal immediately leaves the monitor, *Q* is resumed
- Implemented in other languages including Mesa, C#, Java

Linux Synchronisation

- Linux:
 - Prior to kernel version 2.6, disables interrupts to implement short critical sections
 - Version 2.6 and later, fully preemptive
- Linux provides:
 - atomic integers (set, add, sub, inc, read)
 - mutex locks
 - semaphores
 - spinlocks
 - reader-writer versions of both
- On single-CPU system, spinlocks replaced by enabling and disabling kernel preemption (i.e. disabling interrupts)

Atomic Transactions

- Assures that operations happen as a single logical unit of work, in its entirety, or not at all
- Related to field of database systems
- Challenge is assuring atomicity despite computer system failures
- Must also work with multiple concurrent clients
- Approaches – see textbook for more detail if needed

- Log-based Recovery
- Checkpoints
- Concurrent Atomic Transactions

Deadlocks

- To develop a description of deadlocks, which prevent sets of concurrent processes from completing their tasks
- To present a number of different methods for preventing or avoiding deadlocks in a computer system

System Model

- System consists of resources
- Resource types R_1, R_2, \dots, R_m
 - CPU cycles, memory space, I/O devices
- Each resource type R_i has W_i instances
- Each process utilizes a resource as follows:
 - request
 - use
 - release

Deadlock Characterization

Deadlock can arise if four conditions hold simultaneously.

- **Mutual exclusion:** only one process at a time can use a resource
- **Hold and wait:** a process holding at least one resource is waiting to acquire additional resources held by other processes
- **No preemption:** a resource can be released only voluntarily by the process holding it, after that process has completed its task
- **Circular wait:** there exists a set $\{P_0, P_1, \dots, P_n\}$ of waiting processes such that P_0 is waiting for a resource that is held by P_1 , P_1 is waiting for a resource that is held by P_2 , ..., P_{n-1} is waiting for a resource that is held by P_n , and P_n is waiting for a resource that is held by P_0 .

Deadlock with Mutex Locks

- Deadlocks can occur via system calls, locking, etc

Task 1

```

1 | Lock(Resource_1);
2 | Lock(Resource_2);
3 | // Do Stuff
4 | Unlock(Resource_2);
5 | Unlock(Resource_1);

```

Task 2

```

1 | Lock(Resource_2);
2 | Lock(Resource_1);
3 | // Do Stuff
4 | Unlock(Resource_1);
5 | Unlock(Resource_2);

```

Resource-Allocation Graph

A set of vertices V and a set of edges E

- V is partitioned into two types:
 - $P = \{P_1, P_2, \dots, P_n\}$, the set consisting of all the processes 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 $P_i \rightarrow R_j$
- **Assignment edge** – directed edge $R_j \rightarrow P_i$
- Can be analysed to detect deadlock

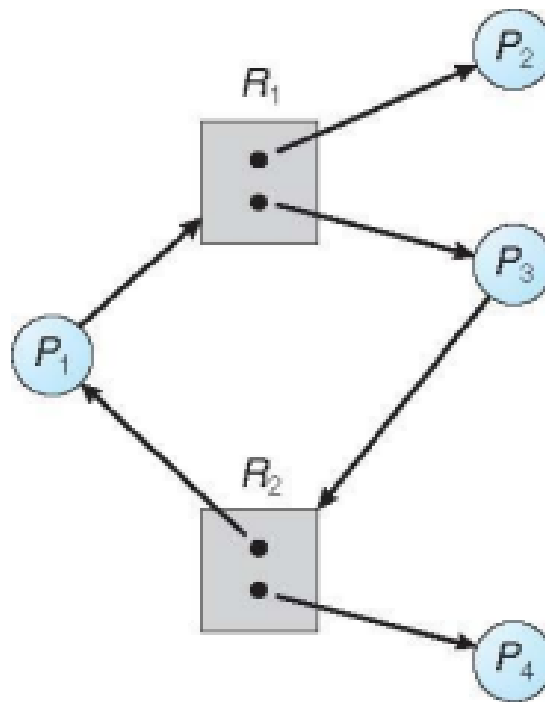


Figure 9: Resource-Allocation Graph

Basic Facts:

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

Methods for Handling Deadlocks

- Ensure that the system will **never** enter a deadlock state
- Allow the system to enter a deadlock state and then recover
- Ignore the problem and pretend that deadlocks never occur in the system used by most operating systems, including UNIX

Deadlock Prevention

Restrain the ways request can be made

- **Mutual Exclusion** – not required for sharable resources; must hold for nonsharable resources
- **Hold and Wait** – must guarantee that whenever a process requests a resource, it does not hold any other resources
 - Require process to request and be allocated all its resources before it begins execution, or allow process to request resources only when the process has none
 - Low resource utilization; starvation possible
- **No Preemption** –
 - If a process that is holding some resources requests another resource that cannot be immediately allocated to it, then all resources currently being held are released
 - Preempted resources are added to the list of resources for which the process is waiting
 - Process will be restarted only when it can regain its old resources, as well as the new ones that it is requesting
- **Circular Wait** – impose a total ordering of all resource types, and require that each process requests resources in an increasing order of enumeration

Deadlock Avoidance

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

- Simplest and most useful model requires that each process declare the **maximum number** of resources of each type that it may need
- The deadlock-avoidance algorithm dynamically examines the resource-allocation state to ensure that there can never be a circular-wait condition
- Resource-allocation *state* is defined by the number of available and allocated resources, and the maximum demands of the processes

Safe State

- When a process requests an available resource, system must decide if immediate allocation leaves the system in a safe state
- System is in **safe state** if there exists a sequence $\langle P_1, P_2, \dots, P_n \rangle$ of ALL processes in the systems such that for each P_i , the resources that P_i can still request can be satisfied by currently available resources + resources held by all the P_j with $j < i$
- That is:
 - If P_i resource needs are not immediately available, then P_i can wait until all P_j have finished
 - When P_j is finished, P_i can obtain needed resources, execute, return allocated resources, and terminate
 - When P_i terminates, P_{i+1} can obtain its needed resources, and so on

Basic Facts:

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

Avoidance Algorithms

- Single instance of a resource type
 - Use a resource-allocation graph
- Multiple instances of a resource type
 - Use the banker's algorithm

Resource-Allocation Graph Scheme

- **Claim edge** $P_i \rightarrow R_j$ indicated that process P_i may request resource R_j ; represented by a dashed line
- Claim edge converts to request edge when a process requests a resource
- Request edge converted to an assignment edge when the resource is allocated to the process
- When a resource is released by a process, assignment edge reconverts to a claim edge
- Resources must be claimed *a priori* in the system
- Suppose that process P_i requests a resource R_j
- The request can be granted only if converting the request edge to an assignment edge does not result in the formation of a cycle in the resource allocation graph

Bankers Algorithm

- Multiple instances
 - Each process must a priori claim maximum use
 - When a process requests a resource it may have to wait
 - When a process gets all its resources it must return them in a finite amount of time
- > Not assessable for exam

Deadlock Detection

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

Single Instance of Each Resource Type

-
- Maintain **wait-for** graph
 - Nodes are processes
 - $P_i \rightarrow P_j$ if P_i is waiting for P_j
 - Periodically invoke an algorithm that searches for a cycle in the graph. 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

Detection-Algorithm Usage

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

Recovery from Deadlock

Process Termination:

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

Resource Preemption:

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

Main Memory

- To provide a detailed description of various ways of organizing memory hardware
- To discuss various memory-management techniques, including paging and segmentation
- To provide a detailed description of the Intel Pentium, which supports both pure segmentation and segmentation with paging

Background

- Program must be brought (from disk) into memory and placed within a process for it to be run
- Main memory and registers are only storage CPU can access directly
- Memory unit only sees a stream of addresses + read requests, or address + data and write requests
- Register access in one CPU clock (or less)
- Main memory can take many cycles, causing a **stall**
- **Cache** sits between main memory and CPU registers
- Protection of memory required to ensure correct operation

Base and Limit Registers

- A pair of **base** and **limit registers** define the logical address space

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

Address Binding

- Programs on disk, ready to be brought into memory to execute form 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?
- Further, 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
 - Linker or loader will bind relocatable addresses to absolute addresses
 - Each binding maps one address space to another

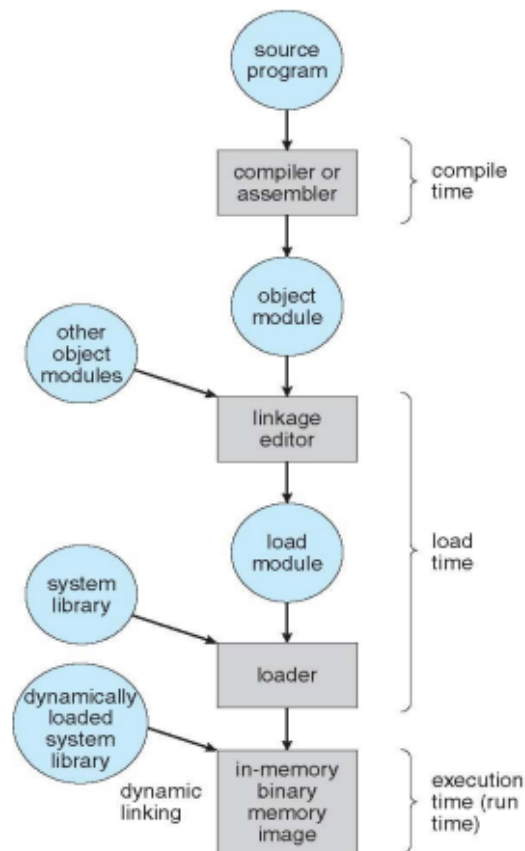


Figure 10: Multistep Processing of a User Program

Logical vs. Physical Address Space

- The concept of a logical address space that is bound to a separate **physical address space** is central to proper memory management
 - **Logical address** – generated by the CPU; also referred to as **virtual address**
 - **Physical address** – address seen by the memory unit
- Logical and physical addresses are the same in compile-time and load-time address-binding schemes; logical (virtual) and physical addresses differ in execution-time address-binding scheme
- **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 the chapter

-
- To start, consider simple scheme where the value in the relocation register is added to every address generated by a user process at the time it is sent to memory
 - Base register now called **relocation register**
 - MS-DOS on Intel 80x86 used 4 relocation registers
 - 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

Dynamic relocation using a relocation register

- 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

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

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
- Does the swapped out process need to swap back in to same physical addresses?
- Depends on address binding method
 - Plus consider pending I/O to/from process memory space
- Modified versions of swapping are found on many systems (i.e. UNIX, Linux, and Windows)
 - Swapping normally disabled
 - Started if more than threshold amount of memory allocated
 - Disabled again once memory demand reduced below threshold

Context Switch Time including Swapping

- If next processes to be put on CPU is not in memory, need to swap out a process and swap in target process
- Context switch time can then be very high

- 100MB process swapping to hard disk with transfer rate of 50MB/sec
 - Swap out time of 2000ms
 - Plus swap in of same sized process
 - Total context switch swapping component time of 4000ms (4 seconds)
- Can reduce if reduce size of memory swapped – by knowing how much memory really being used
 - System calls to inform OS of memory use via `request_memory()` and `release_memory()`
- Other constraints as well on swapping
 - Pending I/O – can't swap out as I/O would occur to wrong process
 - Or always transfer I/O to kernel space, then to I/O device
 - Known as **double buffering**, adds overhead
- Standard swapping not used in modern operating systems
 - But modified version common
 - Swap only when free memory extremely low

Swapping on Mobile Systems

- Not typically supported
 - Flash memory based
 - Small amount of space
 - Limited number of write cycles
 - Poor throughput between flash memory and CPU on mobile platform
- Instead use other methods to free memory if low
 - iOS asks apps to voluntarily relinquish allocated memory
 - Read-only data thrown out and reloaded from flash if needed
 - Failure to free can result in termination
 - Android terminates apps if low free memory, but first writes **application state** to flash for fast restart
 - Both OSes support paging as discussed below

Contiguous Memory Allocation

- 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
- 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 of logical addresses – each logical address must be less than the limit register
 - MMU maps logical address *dynamically*
 - Can then allows actions such as kernel code being **transient** and kernel changing size
- 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:
 - 1) allocated partitions
 - 2) free partitions (hole)

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](#)
- 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

Segmentation

- Memory-management scheme that supports user view of memory
- A program is a collection of segments
 - A segment is a logical unit such as:
 - main program
 - procedure
 - function
 - method
 - object
 - local variables, global variables
 - common block
 - stack
 - symbol table
 - arrays

Segmentation Architecture

- Logical address consists of a two tuples:
 - $\langle \text{segment-number}, \text{offset} \rangle$
- **Segment table** – maps two-dimensional physical addresses; each table entry has:
 - **base** – contains the starting physical address where the segments reside in memory
 - **limit** – specifies the length of the segment
- **Segment-table base register (STBR)** points to the segment table's location in memory
- **Segment-table length register (STLR)** indicates number of segments used by a program;
 - segment number s is legal if $s < STLR$
- Protection
 - With each entry in segment table associate:
 - validation bit = 0 \Rightarrow illegal segment

- read/write/execute privileges
- Protection bits associated with segments; code sharing occurs at segment level
- Since segments vary in length, memory allocation is a dynamic storage-allocation problem

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 addresses 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

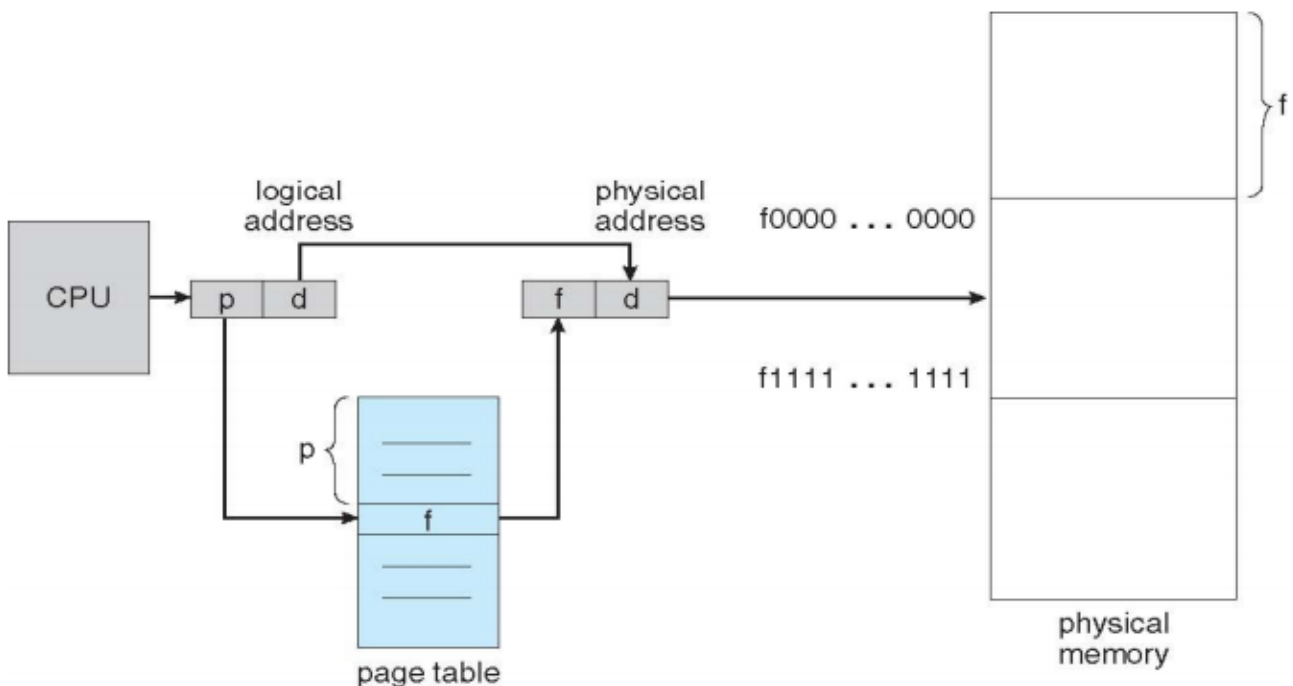


Figure 11: Paging Hardware

- Calculating internal fragmentation
 - 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 size – 8KB and 4MB
- Process view and physical memory now very different
- By implementation process can only access its own memory

Implementation of Page Table

- Page table is kept in main memory
- **Page-table base register (PTBR)** points to the page table
- **Page-table length register (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 **associative memory** or **translation look-aside buffers (TLBs)**
- 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

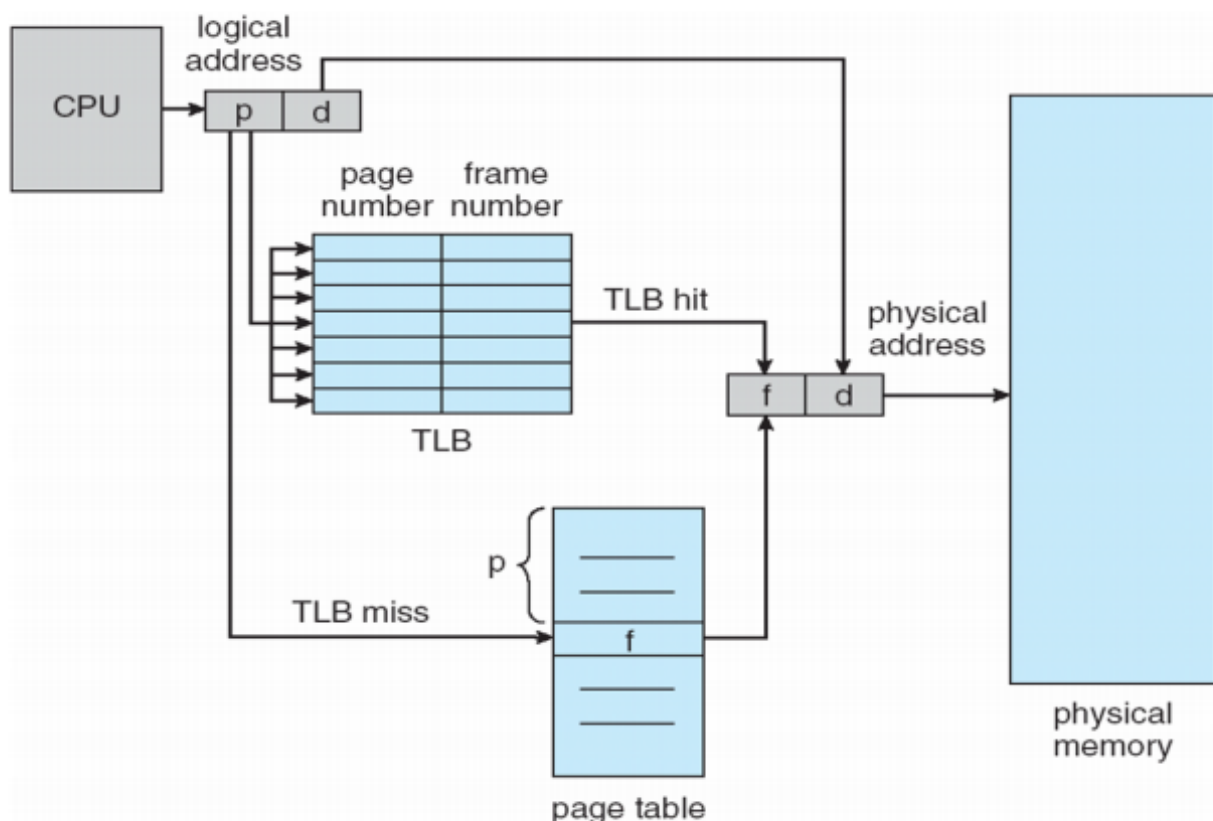


Figure 12: Paging Hardware with TLB

Memory Protection

-
- 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

Shared Pages

- **Shared code**
 - One copy of read-only (**reentrant**) code shared among processes (i.e. text editors, compilers, window systems)
 - Similar to multiple threads sharing the same process space
 - Also useful for interprocess communication if sharing of read-write pages is allowed
- **Private code and data**
 - Each process keeps a separate copy of the code and data
 - The pages for the private code and data can appear anywhere in the logical address space

Structure of the Page Table

- Memory structures for paging can get huge using straight-forward methods
 - Consider a 32-bit logical address space on modern computers
 - Page size of 4KB (2^{12})
 - Page table would have 1 million entries ($2^{32} / 2^{12}$)
 - If each entry is 4 bytes \rightarrow 4MB of physical address space / memory for page table alone
 - That amount of memory used to cost a lot
 - Don't want to allocate that contiguously in main memory

Hierarchical Page Tables

- Break up the logical address space into multiple page tables
- A simple technique is a two-level page table
- We then page the page table

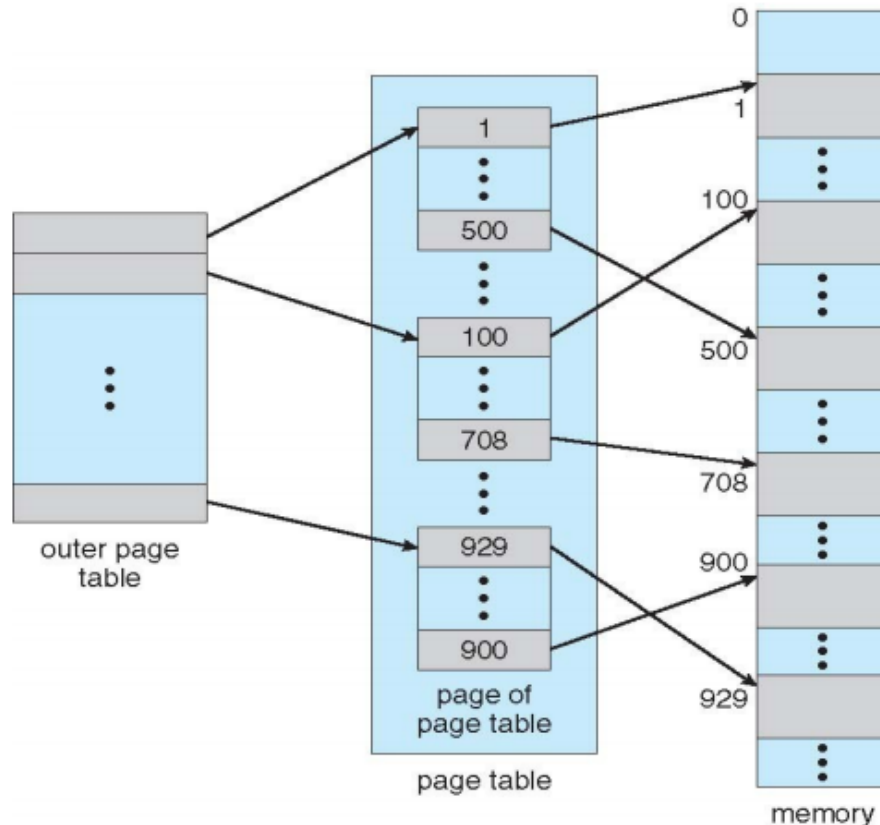


Figure 13: Two-Level Page-Table Scheme

Hashed Page Tables

- Common in address spaces > 32 bits
- The virtual page number is hashed into a page table
 - This page table contains a chain of elements hashing to the same location
- Each element contains
 - 1) The virtual page number
 - 2) The value of the mapped page frame
 - 3) A pointer to the next element
- Virtual page numbers are compared in this chain searching for a match
 - If a match is found, the corresponding physical frame is extracted
- Variation for 64-bit address is **clustered page tables**
 - Similar to hashed but each entry refers to several pages (such as 16) rather than 1
 - Especially useful for **sparse** address spaces (where memory references are non-contiguous and scattered)

Inverted Page Table

- Rather than each process having a page table and keeping track of all possible logical pages, track all physical pages
- One entry for each real page of memory
- Entry consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page
- Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs
- Use hash table to limit the search to one – or at most a few – page-table entries
 - TLB can accelerate access
- But how to implement shared memory?
 - One mapping of a virtual address to the shared physical address

-
- To describe the benefits of a virtual memory system
 - To explain the concepts of demand paging, page-replacement algorithms, and allocation of page frames
 - To discuss the principle of the working-set model
 - To examine the relationship between shared memory and memory-mapped files
 - To explore how kernel memory is managed

Background

- Code needs to be in memory to execute, but entire program rarely used
 - Error code, unusual routines, large data structures
- Entire program code not needed at same time
- Consider ability to execute partially-loaded program
 - Program no longer constrained by limits of physical memory
 - Each program takes less memory while running → more programs run at the same time
 - Increased CPU utilization and throughput with no increase in response time or turnaround time
 - Less I/O needed to load or swap programs into memory → each user program runs faster
- **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 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-address Space

- 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

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 => reference to it

- invalid reference => abort
- not-in-memory => bring to memory
- **Lazy swapper** – never swaps a page into memory unless page will be needed
 - Swapper that deals with pages is a **pager**

Basic Concepts:

- 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 different from non demand-paging
- If page needed and not memory resident
 - Need to detect and load into memory from storage
 - Without changing program behaviour
 - Without programmer needing to change code

Page Fault

- If there is a reference to a page, first reference to that page will trap to operating system: **page fault**

1) Operating system looks at another table to decide:

- Invalid reference => abort
- Just not in memory

1) Find free frame

2) Swap page into frame via scheduled disk operation

3) Reset tables to indicate page now in memory. Set validation bit = v

4) Restart the instruction that caused the page fault

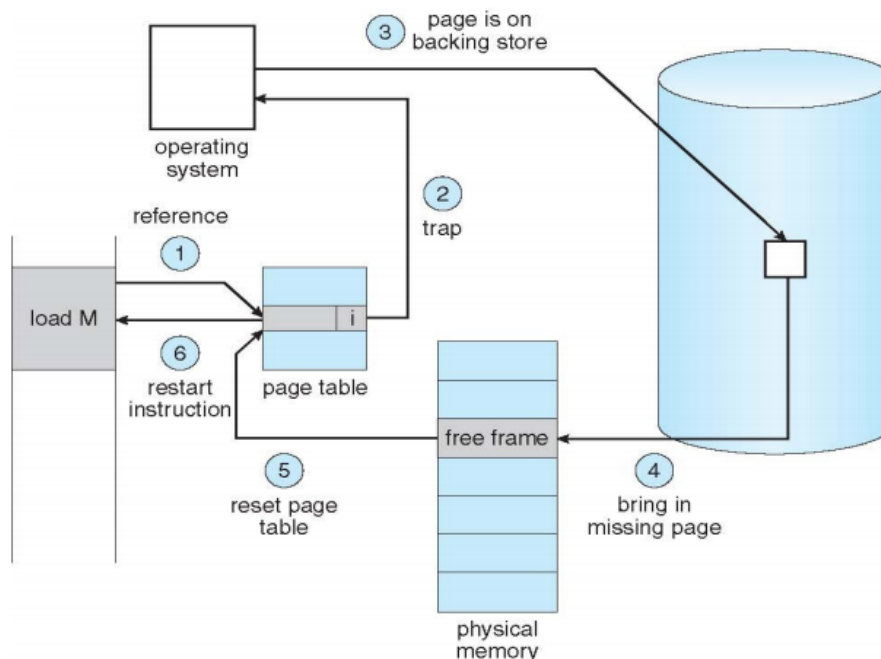


Figure 14: Steps in Handling a Page Fault

Aspects of Demand Paging

- 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

Performance of Demand Paging

- Stages in Demand Paging (worse case)
 - 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:
 - 1) Wait in a queue for this device until the read request is serviced
 - 2) Wait for the device seek and/or latency time
 - 3) Begin the transfer of the page to a free frame
 - 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
- 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)

> $EAT = (1 - p) \times \text{memory access} + p (\text{page fault overhead} + \text{swap page out} + \text{swap page in})$

Demand Paging Optimizations

- Swap space I/O faster than file system I/O even if on the same device
 - Swap allocated in larger chunks, less management needed than file system
- Copy entire process image to swap space at process load time
 - Then page in and out of swap space
 - Used in older BSD Unix
- Demand page in from program binary on disk, but discard rather than paging out when freeing frame
 - Used in Solaris and current BSD
 - Still needed to write to swap space
 - Pages not associated with a file (like stack and heap) – **anonymous memory**
 - Pages modified in memory but not yet written back to the file system
- Mobile systems
 - Typically don't support swapping
 - Instead, demand page from file system and reclaim read-only pages (such as code)

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

What happens if there is no free frame?

- Used up by process pages
- Also in demand from the kernel, I/O buffers, etc
- How much to allocate to each?
- Page replacement – find some page in memory, but not really in use, page it out
 - Algorithm – terminate? swap out? replace the page?
 - Performance – want an algorithm which will result in minimum number of page faults
- Same page may be brought into memory several times

Page Replacement

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

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
 - 1) Bring the desired page into the (newly) free frame; update the page and frame tables
 - 2) Continue the process by restarting the instruction that caused the trap
- > Note now potentially 2 page transfers for page fault – increasing EAT

Page and Frame Replacement Algorithms

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

First-In-First-Out (FIFO) Algorithm

- Adding more frames can cause more page faults!
 - **Beladys Anomaly**
- How to track ages of pages?
 - Just use a FIFO queue

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

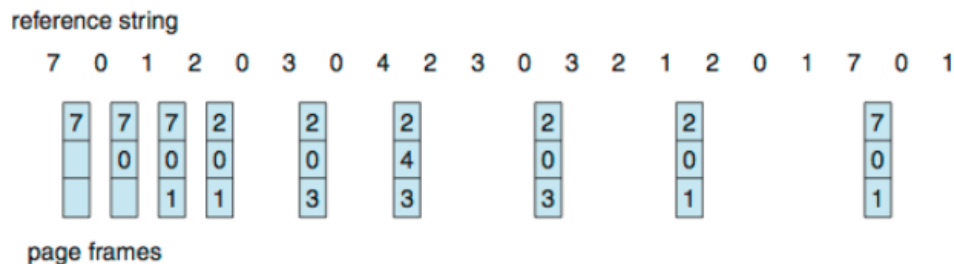


Figure 15: Optimal

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
- Generally good algorithm and frequently used

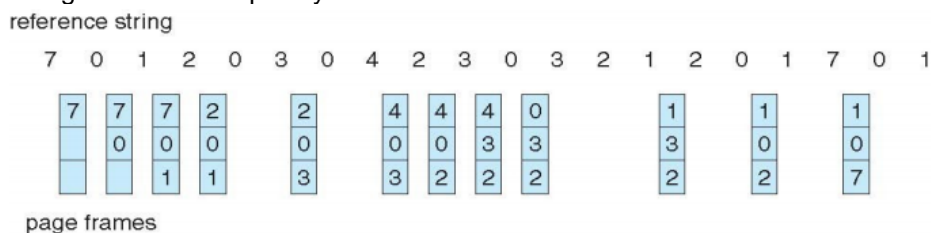


Figure 16: LRU

- 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 and OPT are cases of **stack algorithms** that don't have Beladys Anomaly

LRU Approximation Algorithms

- LRU needs special hardware and still slow
- **Reference bit**
 - With each page associate a bit, initially = 0

- When page is referenced bit set to 1
- Replace any with reference bit = 0 (if one exists)
 - We do not know the order, however
- **Second-chance algorithm**
 - Generally FIFO, plus hardware-provided reference bit
 - **Clock** replacement
 - If page to be replaced has
 - Reference bit = 0 → replace it
 - Reference bit = 1 then:
 - Set reference bit 0, leave page in memory
 - Replace next page, subject to same rules

Enhanced Second-Chance Algorithm

- Improve algorithm by using reference bit and modify bit (if available) in concert
- Take ordered pair (reference, modify)
- 1) (0, 0) neither recently used nor modified – best page to replace
- 2) (0, 1) not recently used but modified – not quite as good, must write out before replacement
- 3) (1, 0) recently used but clean – probably will be used again soon
- 4) (1, 1) recently used and modified – probably will be used again soon and need to write out before replacement
- When page replacement called for, use the clock scheme but use the four classes replace page in lowest non-empty class
 - Might need to search circular queue several times

Counting Algorithms

- Keep a counter of the number of references that have been made to each page
 - Not common
- **Least Frequently Used (LFU) Algorithm**: replaces page with smallest count
- **Most Frequently Used (MFU) Algorithm**: based on the argument that the page with the smallest count was probably just brought in and has yet to be used

Page-Buffering Algorithms

- Keep a pool of free frames, always
 - Then frame available when needed; not found at fault time
 - Read page into free frame and select victim to evict and add to free pool
 - When convenient, evict victim
- Possibly, keep list of modified pages
 - When backing store otherwise idle, write pages there and set to non-dirty
- Possibly, keep free frame contents intact and note what is in them
 - If referenced again before reused, no need to load contents again from disk
 - Generally useful to reduce penalty if wrong victim frame selected

Non-Uniform Memory Access

- So far all memory accessed equally
- Many systems are **NUMA** – speed of access to memory varies
 - Consider system boards containing CPUs and memory, interconnected over a system bus
- Optimal performance comes from allocating memory "close to" the CPU on which the thread is scheduled
 - And modifying the scheduler to schedule the thread on the same system board when possible
- Solved by Solaris by creating **lggroups**
 - Structure to track CPU/Memory low latency groups
 - Used by scheduler and pager

- When possible schedule all threads of a process and allocate all memory for that process within the lgroup

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** = a process is busy swapping pages in and out

Demand Paging and Thrashing

- Why does demand paging work? [Locality model](#)
 - Process migrates from one locality to another
 - Localities may overlap
- Why does thrashing occur?
 - $\sigma \text{size of locality} > \text{total memory size}$
 - Limit effects by using local or priority page replacement

Working-Set Model

- Δ = working-set window = a fixed number of page references
- WSS_i (working set of Process P_i) = total number of pages referenced in the most recent Δ (varies in time)
 - if Δ too small will not encompass entire locality
 - if Δ too large will encompass several localities
 - if $\Delta = \infty \Rightarrow$ will encompass entire program
- $D = \sum WSS_i$ = total demand frames
 - Approximation of locality
- if $D > m \Rightarrow$ Thrashing
- Policy if $D > m$, then suspend or swap out one of the processes

page reference table

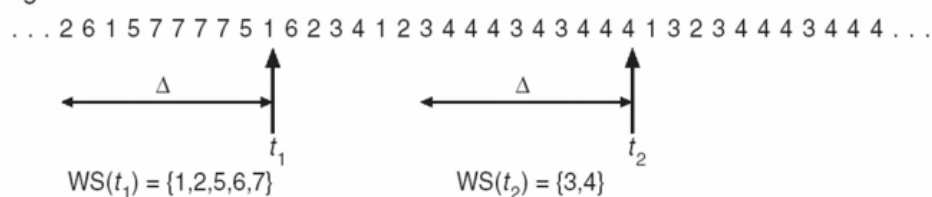


Figure 17: Working-Set Model

Memory-Mapped Files

- Memory-mapped file I/O allows file I/O to be treated as routine memory access by [mapping](#) a disk block to a page in memory
- A file is initially read using demand paging
 - A page-size portion of the file is read from the file system into a physical page
 - Subsequent reads/writes to/from the file are treated as ordinary memory accesses
- Simplifies and speeds file access by driving file I/O through memory rather than `read()` and `write()` system calls
- Also allows several processes to map the same file allowing the pages in memory to be shared
- But when does written data make it to disk?
 - Periodically and/or at file `close()` time

-
- For example, when the pager scans for dirty pages

Allocating Kernel Memory

- Treated differently from user memory
- Often allocated from a free-memory pool
 - Kernel requests memory for structures of varying sizes
 - Some kernel memory needs to be contiguous (i.e. for device I/O)

Buddy System

- Allocates memory from fixed-size segment consisting of physically-contiguous pages
- Memory allocated using power-of-2 allocator
 - Satisfies requests in units sized as power of 2
 - Request rounded up to next highest power of 2
 - When smaller allocation needed than is available, current chunk split into two buddies of next-lower power of 2
 - Continue until appropriate sized chunk available
- For example, assume 256KB chunk available, kernel requests 21KB
 - Split into A_L and A_R of 128KB each
 - One further divided into B_L and B_R of 64KB
 - One further into C_L and C_R of 32KB each – one used to satisfy request
- Advantage – quickly **coalesce** unused chunks into larger chunk
- Disadvantage – fragmentation

Slab Allocator

- Alternate strategy
- **Slab** is one or more physically contiguous pages
- **Cache** consists of one or more slabs
- Single cache for each unique kernel data structure
 - Each cache filled with **objects** – instantiations of the data structure
- When cache created, filled with objects marked as **free**
- When structures stored, objects marked as **used**
- If slab is full of used objects, next object allocated from empty slab
 - If no empty slabs, new slab allocated
- Benefits include no fragmentation, fast memory request satisfaction

Other Issues

Page Size

- Sometimes OS designers have a choice
 - Especially if running on custom-built CPU
- Page size selection must take into consideration:
 - Fragmentation
 - Page table size
 - **Resolution**
 - I/O overhead
 - Number of page faults
 - Locality
 - TLB size and effectiveness
- Always power of 2, usually in the range 2^{12} (4,096 bytes) to 2^{22} (4,194,304 bytes)
- On average, growing over time

TLB Reach

- TLB Reach – the amount of memory accessible from the TLB
- $\text{TLB Reach} = (\text{TLB Size}) \times (\text{Page Size})$
- Ideally, the working set of each process is stored in the TLB
 - Otherwise there is a high degree of page faults
- Increase the Page Size
 - This may lead to an increase in fragmentation as not all applications require a large page size
- Provide Multiple Page Sizes
 - This allows applications that require larger page sizes the opportunity to use them without an increase in fragmentation

Program Structure

- Program structure
 - `int[128,128] data;`
 - Each row is stored in one page
- Program 1

```
1 | for (j = 0; j < 128; j++) {
2 |     for (i = 0; i < 128; i++) {
3 |         data[i,j] = 0;
4 |     }
5 | }
```

 - $128 \times 128 = 16,384$ page faults
- Program 2

```
1 | for (i = 0; i < 128; i++) {
2 |     for (j = 0; j < 128; j++) {
3 |         data[i,j] = 0;
4 |     }
5 | }
```

 - 128 page faults

I/O interlock

- **I/O interlock** – Pages must sometimes be locked into memory
- Consider I/O – Pages that are used for copying a file from a device must be locked from being selected for eviction by a page replacement algorithm
- **Pinning** of pages to lock into memory

Mass-Storage System

- To describe the physical structure of secondary storage devices and its effects on the uses of the devices
- To explain the performance characteristics of mass-storage devices
- To evaluate disk scheduling algorithms
- To discuss operating-system services provided for mass storage, including RAID

Overview of Mass Storage Structure

- **Magnetic disks** provide bulk of secondary storage of modern computers
 - Drives rotate at 60 to 250 times per second
 - **Transfer rate** is rate at which data flow between drive and computer
 - **Positioning time** (**random-access time**) is time to move disk arm to desired cylinder (**seek time**) and time for desired sector to rotate under the disk head (**rotational latency**)
 - **Head crash** results from disk head making contact with the disk surface
 - That's bad
- Disks can be removable
- Drive attached to computer via **I/O bus**

-
- Busses vary, including EIDE, ATA, SATA, USB, Fibre Channel, SCSI, SAS, Firewire
 - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array

Magnetic Disks

- Platters range from .85" to 14" (historically)
 - Commonly 3.5", 2.5" and 1.8"
- Range from 30GB to 3TB per drive
- Performance
 - Transfer rate: 6Gb/s
 - Effective transfer: 1Gb/s
 - Seek time 3ms to 12ms (9ms common)
 - Average seek time measured or calculated based on 1/3 of tracks
 - Latency based on spindle speed
 - $1 / (\text{RPM} / 60) = 60 / \text{RPM}$
 - Average latency = 1/2 latency

Solid-State Disks

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

SSHD

- Solid State Hybrid Drives, e.g. Apple Fusion drive
- Incorporate Flash Drive and Magnetic Disk Drive
- May use Flash as a cache, or else may be two separate drives with one interface
- Reduced power consumption for laptops
- May be manual or automatic configuration (e.g. OS moves commonly used files to SSD)

Magnetic Tape

- Was early secondary-storage medium
 - Evolved from open spools to cartridges
- Relatively permanent and holds large quantities of data
- Access time slow
- Random access ~1000 times slower than disk
- Mainly used for backup, storage of infrequently-used data, transfer medium between systems
- Kept in spool and wound or rewound past read-write head
- Once data under head, transfer rates comparable to disk
 - 140MB/s and greater
- 200GB to 1.5TB typical storage
- Robot tape drives can store and access thousands of tapes, with multi-second latency

Disk Structure

- Disk drives are addressed as large 1-dimensional arrays of **logical blocks**, where the logical block is the smallest unit of transfer

- Low-level formatting creates **logical blocks** on physical media
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - Sector 0 is the first sector of the first track on the outermost cylinder
 - Mapping proceeds in order through that track, then the rest of the tracks in the cylinder, and then through the rest of the cylinders from outermost to innermost
- Logical to physical address should be easy
 - Except for bad sectors
 - Non-constant number of sectors per track via constant angular velocity
- Older disks would allow physical addressing (Cylinder Number, Sector Number, Head Number), but most disks now only have (linear) logical addressing. This means the OS does not know exactly where a disk block is located

Disk Attachment

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

Storage Array

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

Storage Area Network

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

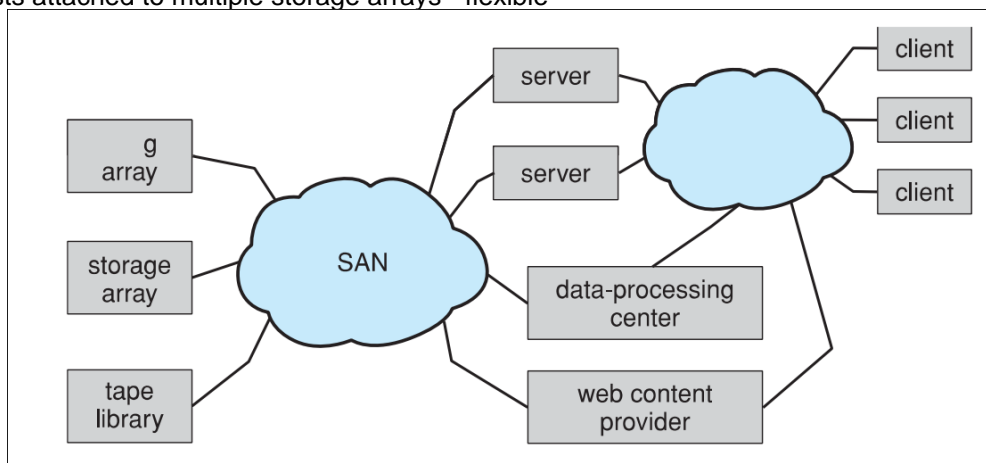


Figure 18: Storage Area Network

- SAN is one or more storage arrays
 - Connected to one or more Fibre Channel switches

-
- Hosts also attach to the switches
 - Storage made available via **LUN (Logical Unit Number) Masking** from specific arrays to specific servers
 - Easy to add or remove storage, add new host and allocate it storage
 - Over low-latency Fibre Channel fabric

Network-Attached Storage

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

Disk Scheduling

- The operating system is responsible for using hardware efficiently – for the disk drives, this means having a fast access time and disk bandwidth
- Minimize seek time
- Seek time = seek distance
- Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer
- There are many sources of disk I/O request
 - OS
 - System processes
 - Users processes
- I/O request includes input or output mode, disk address, number of sectors to transfer
- OS maintains queue of requests, per disk or device
- Idle disk can immediately work on I/O request, busy disk means work must queue
 - Optimization algorithms only make sense when a queue exists
- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying "depth")
- Several algorithms exist to schedule the servicing of disk I/O requests
- The analysis is true for one or many platters
- We illustrate scheduling algorithms with a request queue (0-199)

FCFS

Obvious

SSTF

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

SCAN

- The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues
- **SCAN algorithm** sometimes called the **elevator algorithm**

C-SCAN

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

C-Look

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

Selecting a Disk-Scheduling Algorithm

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

Disk Management

- **Low-level formatting** or **physcial formatting** – Dividing a disk into sectors that the disk controller can read or write
 - Each sector can hold header information, plus data, plus error correction code (**ECC**)
 - Usually 512 bytes of data but can be selectable
- To use a disk to hold files, the operating system still needs to record its own data structures on the disk
 - **Partition** the disk into one or more groups of cylinders, each treated as a logical disk
 - **Logical formatting** or "making a file system"
 - To increase efficiency most file systems group blocks into **clusters**
 - Disk I/O done in blocks
 - File I/O done in clusters
- Raw disk access for apps that want to do their own block management, keep OS out of the way (databases for example)
- Boot block initializes system
 - The bootstrap is stored in ROM
 - **Bootstrap loader** program stored in boot blocks of boot partition
- Methods such as **sector sparing** used to handle bad blocks

Swap-Space Management

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

-
- What if a system runs out of swap space?
 - Some systems allow multiple swap spaces

RAID Structure

- RAID – redundant array of inexpensive disks
 - multiple disk drives provides reliability via **redundancy**
- Increases the **mean time to failure**
- **Mean time to repair** – exposure time when another failure could cause data loss
- **Mean time to data loss** based on above factors
- If mirrored disks fail independently, consider disk with 1,300,000 mean time to failure and 10 hour mean time to repair
 - Mean time to data loss is $100,000^2 / (2 \times 10)$, 57,000 years
- Frequently combined with **NVRAM** to improve write performance
- RAID is arranged into six different levels
- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively
- Disk **striping** uses a group of disks as one storage unit
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
 - **Mirror** or **shadowing** (**RAID 1**) keeps duplicate of each disk
 - Striped mirrors (**RAID 1+0**) or mirrored stripes (**RAID 0+1**) provides high performance and high reliability
 - **Block interleaved parity** (**RAID, 4, 5, 6**) uses much less redundancy
- RAID within a storage array can still fail if the array fails, so automatic **replication** of the data between arrays is common
- Frequently, a small number of **hot-spare** disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them

Other Features

- Regardless of where RAID implemented, other useful features can be added
- **Snapshot** is a view of the file system before a set of changes take place (i.e. at a point in time)
- Replication is automatic duplication of writes between separate sites
 - For redundancy and disaster recovery
 - Can be synchronous or asynchronous
- Hot spare disk is unused, automatically used by RAID production if a disk fails to replace the failed disk and rebuild the RAID set if possible
 - Decreases mean time to repair

Extensions

- RAID alone does not prevent or detect data corruption or other errors, just disk failures
- Solaris ZFS adds **checksums** of all data and metadata
- Checksums kept with pointer to object, to detect if object is the right one and whether it changed
- Can detect and correct data and metadata corruption
- ZFS also removes volumes, partitions
 - Disks allocated in **pools**
 - Filesystems with a pool share that pool, use and release space like **malloc()** and **free()** memory allocate / release calls

Stable-Storage Implementation

- Write-ahead log scheme requires stable storage
- Stable storage means data is never lost (due to failure, etc)
- To implement stable storage:
 - Replicate information on more than one nonvolatile storage media with independent failure modes

-
- Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery
 - Disk write has 1 of 3 outcomes
 - 1) **Successful completion** – The data were written correctly on disk
 - 2) **Partial failure** – A failure occurred in the midst of transfer, so only some of the sectors were written with the new data, and the sector being written during the failure may have been corrupted
 - 3) **Total failure** – The failure occurred before the disk write started, so the previous data values on the disk remain intact
 - If failure occurs during block write, recovery procedure restores block to consistent state
 - System maintains 2 physical blocks per logical block and does the following:
 - 1) Write to 1st physical
 - 2) When successful, write to 2nd physical
 - 3) Declare complete only after second write completes successfully
 - System frequently use NVRAM as one physical to accelerate

File-System Interface

- To explain the function of file systems
- To describe the interfaces to file systems
- To discuss file-system design tradeoffs, including access methods, file sharing, file locking, and directory structures
- To explore file-system protection

File Concept

- Contiguous logical address space
- Types:
 - Data
 - numeric
 - character
 - binary
 - Program
- Contents defined by file's creator
 - Many types
 - Consider [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, data, and user identification** – data for protection, security, and usage monitoring
- Information about files are kept in the directory structure, which is maintained on the disk
- Many variations, including extended file attributes such as file checksum
- Information kept in the directory structure

File Operations

- File is an [abstract data type](#)
- [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

Close

- 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

Open File Locking

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

File Structure

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

Access Methods

- Sequential Access

```

1 | read next
2 | write next
3 | reset

```

no read after last write (rewrite)

- Direct Access – file is fixed length **logical records**

```

1 | read n
2 | write n
3 | position to n
4 |     read next
5 |     write next
6 | rewrite n
n = ~( )relative block number

```

- Relative block numbers allow OS to decide where file should be placed

Disk and Directory Structure

Directory Structure

- A collection of nodes containing information about all files
- > Both directory structure and files reside on disk

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 known as a **volume**
- Each volume containing file system also tracks that file system's info in **device directory** or **volume table of contents**
- As well as **general-purpose file systems** there are many **special-purpose file systems**, frequently all within the same operating system or computer

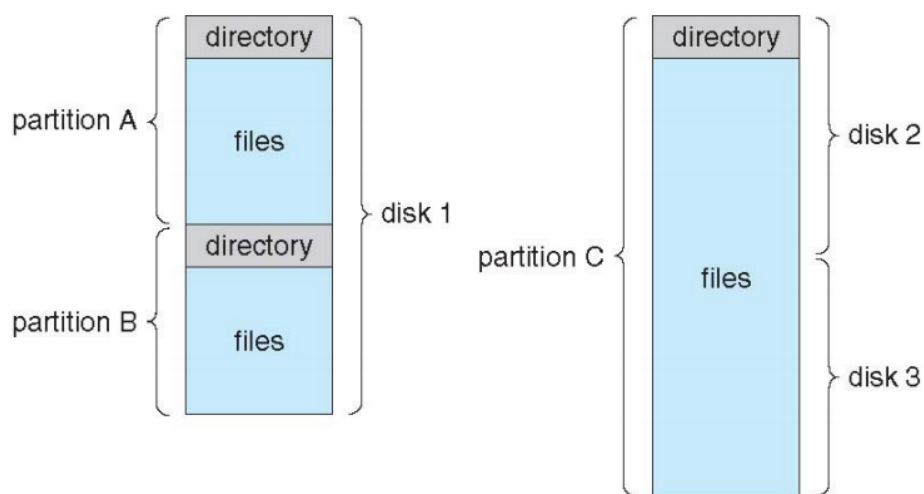


Figure 19: A Typical File-system Organization

Operations Performed on Directory

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

Organize the Directory (Logically) to Obtain

- Efficiency – locating a file quickly
- Naming – convenient to users

-
- Two users can have same name for different files
 - The same file can have several different names
 - Grouping – logical grouping of files by properties, (e.g. all Java programs, all games, ...)

File Sharing

- Sharing of files on multi-user systems is desirable
- Sharing may be done through a **protection** scheme
- On distributed systems, files may be shared across a network
- Network File System (NFS) is a common distributed file-sharing method
- If multi-user system
 - **User IDs** identify users, allowing permissions and protections to be per-user
 - **Group IDs** allow users to be in groups, permitting group access rights
 - Owner of a file / directory
 - Group of a file / directory

Remote File Systems

- Uses networking to allow file system access between systems
 - Manually via programs like FTP
 - Automatically, seamlessly using **distributed file systems**
 - Semi automatically via the **world wide web**
- **Client-server** model allows clients to mount remote file systems from servers
 - Server can serve multiple clients
 - Client and user-on-client identification is insecure or complicated
 - **NFS** is standard UNIX client-server file sharing protocol
 - **CIFS** is standard Windows protocol
 - Standard operating system file calls are translated into remote calls
- Distributed Information Systems (**distributed naming services**) such as LDAP, DNS, NIS, Active Directory implement unified access to information needed for remote computing

Failure Modes

- All file systems have failure modes
 - For example corruption of directory structures or other non-user data, called **metadata**
- Remote file systems add new failure modes, due to network failure, server failure
- Recovery from failure can involve **state information** about status of each remote request
- **Stateless** protocols such as NFS v3 include all information in each request, allowing easy recovery but less security

Protection

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

Access Lists and Groups

- Mode of access: read, write, execute
- Three classes of users on Unix/Linux
 - 1) owner access – 7
 - 2) group access – 6
 - 3) public access – 1
- Ask manager to create a group (unique name), say G, and add some users to the group
- For a particular file (say game) or subdirectory, define an appropriate access

File System Implementation

- To describe the details of implementing local file systems and directory structures
- To describe the implementation of remote file systems
- To discuss block allocation and free-block algorithms and trade-offs

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** – storage structure consisting of information about a file
- **Device driver** controls the physical device
- File system organized into layers

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 number to physical block number
 - Manages free space, disk allocation
- **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
- 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 40 types, with **extended file system** ext2 and ext3 leading; plus distributed file systems, etc)
 - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE

File-System Implementation

- 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 number of blocks, number of free blocks, block size, free block pointers or array
- Directory structure organizes the files
 - Names and inode numbers, file file table
- Per-file **File control block** ([FCB](#)) contains many details about the file
 - inode number, permissions, size dates

In-Memory File System Structures

- Mount table storing file system mounts, mount points, file system types
- The following figure illustrates the necessary file system structures provided by the operating systems
- Figure a refers to opening a file
- Figure b refers to reading a file
- Plus buffers hold data blocks from secondary storage
- Open returns a file handle for subsequent use
- Data from read eventually copied to specific user process memory address

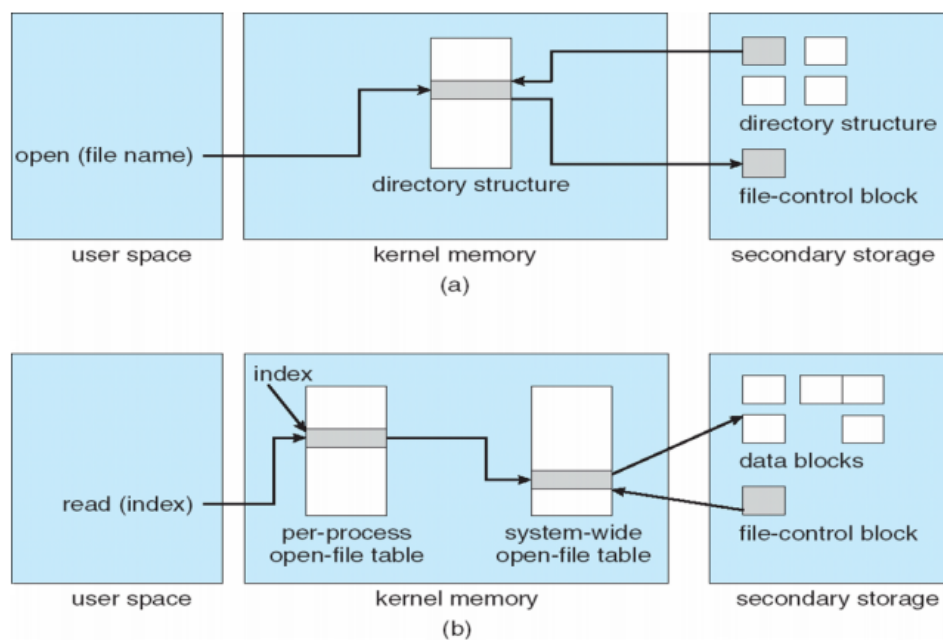


Figure 20: In-Memory File System Structures

Directory Implementation

- **Linear list** of file names with pointer to the data blocks
 - Simple to program
 - Time-consuming to execute
 - Linear search time
 - Could keep ordered alphabetically via linked list or use B+ tree
- **Hash Table** – linear list with hash data structure
 - Decreases directory search time
 - [Collisions](#) – situations where two file names hash to the same location
 - Only good if entries are fixed size, or use chained-overflow method

Allocation Methods

Contiguous

- An allocation method refers to how disk blocks are allocated for file
- **Contiguous allocation** – each file occupies set of contiguous blocks
 - Best performance in most cases
 - Simple – only starting location (block number) and length (number of blocks) are required
 - Problems include finding space for file, knowing file size, external fragmentation, need for [compaction off-line \(downtime\)](#) or [on-line](#)

Extent-Based Systems

- Many newer file systems (i.e. Veritas File System) use a modified contiguous allocation scheme
- Extent-based file systems allocate disk blocks in extents
- An **extent** is a contiguous blocks of disks
 - Extents are allocated for file allocation
 - A file consists of one or more extents

Linked

- **Linked allocation** – each file 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
- FAT (File Allocation Table) variation
 - Beginning of volume has table, indexed by block number
 - Much like a linked list, but faster on disk and cacheable
 - New block allocation simple

Indexed

- **Indexed allocation** – each file has its own [index block\(s\)](#) for pointers to its data block
- Need index table
- Random access
- Dynamic access without external fragmentation, but have overhead of index block
- Mapping from logical to physical in a file of maximum size of 256K bytes and block size of 512 bytes. We need only 1 block for index table

Performance

- Best method depends on file access type
 - Contiguous great for sequential and random
- Linked good for sequential, not random
- Declare access type at creation → select either contiguous or linked
- Indexed more complex
 - Single block access could require 2 index block reads then data block read
 - Clustering can help improve throughput, reduce CPU overhead

Free-Space Management

- File system maintains **free-space list** to track available blocks/clusters
 - (Using term "block" for simplicity)

-
- **Bit vector** or **bit map** (n blocks)
 - Grouping
 - Modify linked list to store address of next $n-1$ free blocks in first free block, plus a pointer to next block that contains free-block-pointers (like this one)
 - Counting
 - Because space is frequently contiguously used and freed, with contiguous-allocation, extents, or clustering
 - Keep address of first free block and count of following free blocks
 - Free space list then has entries containing addresses and counts

Efficiency and Performance

- Efficiency dependent on:
 - Disk allocation and directory algorithms
 - Types of data kept in file's directory entry
 - Pre-allocation or as-needed allocation of metadata structures
 - Fixed-size or varying-size data structures
- Performance
 - Keeping data and metadata close together
 - **Buffer cache** – separate section of main memory for frequently used blocks
 - **Synchronous** writes sometimes requested by apps or needed by OS
 - No buffering / caching – writes must hit disk before acknowledgement
 - **Asynchronous** writes more common, buffer-able, faster
 - **Free-behind** and **read-ahead** – techniques to optimize sequential access
 - Reads frequently slower than writes

Recovery

- **Consistency checking** – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies
 - Can be slow and sometimes fails
- Use system programs to **back up** data from disk to another storage device (magnetic tape, other magnetic disk, optical)
- Recover lost file or disk by **restoring** data from backup

Log Structured File Systems

- **Log structured** (or **journaling**) file systems record each metadata update to the file system as a **transaction**
- All transactions are written to a log
 - A transaction is considered committed once it is written to the log (sequentially)
 - Sometimes to a separate device or section of disk
 - However, the file system may not yet be updated
- The transactions in the log are asynchronously written to the file system structures
 - When the file system structures are modified, the transaction is removed from the log
- If the file system crashes, all remaining transactions in the log must still be performed
- Faster recovery from crash, removes chance of inconsistency of metadata

NFS

- An implementation and a specification of a software system for accessing remote files across LANs (or WANs)
- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP protocol and Ethernet)
- NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures; the NFS specifications independent of these media
- The independence is achieved through the use of RPIC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces

-
- The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services
 - NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures; the NFS specifications independent of these media
 - This independence is achieved through the use of RPC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces
 - The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services

NFS Protocol

- Provides a set of remote procedure calls for remote file operations. The procedures support the following operations:
 - searching for a file within a directory
 - reading a set of directory entries
 - manipulating links and directories
 - accessing file attributes
 - reading and writing files
- NFS servers are [stateless](#); each request has to provide a full set of arguments (NFS V4 is just coming available – very different, stateful)
- Modified data must be committed to the server's disk before results are returned to the client (lose advantages of caching)
- The NFS protocol does not provide concurrency-control mechanisms