

Module I
Unit 1 j Overview of Operating System 1-1 to 1-30

Syllabus * Introduction ; Operating System Structure and **operations, Process management, Memory management, storage** management, Protection and security, Distributed and special purpose Systems; System Structure: Operating system services and interface, System calls and its types. System programs, Operating System Design and implementation, OS structure. Virtual machines, OS debugging and generation, System boot.

✓	Syllabus Topic : Introduction	1-1
1.1	Introduction to Operating System	1-1
i	(Dec. 14, June 15, Nov. 15)	1-1
z	Syllabus Topic : Operating System Structure	1-1
1.2	Operating System Structure	1-1
12. 1	Monolithic Systems	1-2
X22	Layered Systems.	1-2
1.2.3	Client-Server Model.....	1-3
1.2.4	Monolithic Kernel vs. Microkernel	1-4
z	Syllabus Topic : Operating System Operations	1-5
1.3	Operating System Operations	1-5
1.3.1	Dual Mode Operation	1-5
1.3.2	Timer	1-6
z	Syllabus Topic : Process Management	1-6
1.4	Process Management	1-6
✓	Syllabus Topic : Memory Management	V7
1.5	Memory Management	1-7
z	Syllabus Topic : Storage Management	1-7
1.8	Storage Management	1-7
1.6.1	File Management	1-0
1.6.2	Mass Storage Management	1-8
1.6.3	Caching	1-8
1.6.4	M3 Systems	1-9
z	Syliatxis Topic : Protection and Security	1-9
1.7	Protection and Security	1-9
1.7.1	Threats	1-10
z	Syllabus Topic : Distributed System	1-11
1.8	Distributed System	1-11
✓	Syllabus Topic : Special Purpose Systems	1-13
1.9	Special Purpose Systems	1-13
1.9.1	Real-Time Embedded Systems	1-13
1.9.2	Multimedia Systems	1-14
1.9.3	Kan d Systems	1-14
✓	Syllabus Topic : Operating System Services	1-15

1.10	Operating System Services (May 16)	1-15
1.11	Objectives of Operating System	1-16
(June 15, Nov. 15)	u .. b+n+ a.....	1-16
1.11.1	Functions of Operating System	1-16
(Dec. 14, June 15, Nov. 15)1-16	1-16
z	Syllabus Topic : Operating-System Interface	1-17
1.12	User Operating- System Interface	1-17
z*	Syllabus Topic : System Calls	1-18
1.13	System Calls (June 15, Nov. 15, May 16)	vie
z -	Syllabus Topic : Types of System Calls	1-19
1.13.1	Types of System Calls (Juno 15, Nov. 15)	1-19
1.13.2	Some Examples of System Calls —	1-20
z	Syllabus Topic : System Programs	1-20
1.14	System Programs	1-20
1.14.1	Comparison between System Program and Application Program	1-21
z	Syllabus Topic : Operating System Design and Implementation	1-22
1.15	Operating System Design and Implementation	1-22
1.15.1	Design Goals	1-22
1.15.2	Separating Policies from Mechanisms	1-22
1.15.3	Implementation	1-22
z	Syllabus Topic : Virtual Machines	1-23
1.16	Virtual Machines	1-23
1.16.1	History	1-23
1.16.2	Benefits	1-24
1.16.3	Simulation	1-25
1.16.4	Para-Virtualization.	1-25
1.16.5	Implementation	1-25
1.16.6	Examples	1-25
1.16.6(A)	VMware	1-25
1.16.6(B)	The Java Virtual Machine	1-26
1.16.6(C)	The .NET Framework...	1-26
z	Syllabus Topic : Operating System Debugging	1-27
1.17	Operating System Debugging..	1-27
1.17.1	Failure Analysis	1-27
1.17.2	Performance Timing	V27
1.17.3	DTrace.....	1-27
z	Syllabus Topic : Operating System Generation	1-28
1.18	Operating System Generation	v28
z	Syllabus Topic : System Boot	1-28
1.19	System Boot	1-28
1.20	Exam Pack (University and Review Questions)	1-29

Module II

2-1 to 2-32

Chapter 2 : Process Management

Syllabus	: Process concept : Process Scheduling, Operations on process and Interprocess communication; Multithreading.	10
	: Multithreading models and thread libraries, threading issues.	11
Process Scheduling:	Basic concepts. Scheduling algorithms and Criteria. Thread Scheduling and Multiple Processor Scheduling, —	12
2.1	Introduction	2*1
z	Syllabus Topic : Process Concept	2-1
2.2	Process Concept (Dec. 14)	2-1
2.3	Context Switch	2-2
✓	Syllabus Topic : Operation on Processes	2-2
2.4	Operations on Processes	2-2
2.4.1	Process Creation	2-2
2.4.2	Process Termination./.	2-3
2.5	Process Control Block (Dec. 14)	2-3
2.6	Process States and Process State Transition Diagram	2-4
2.7	Process vs. Thread	2-5
✓	Syllabus Topic : Process Scheduling	2-6
2.8	Process Scheduling	2-6
2.8.1	Scheduling Queues and Schedulers	2-6
2.8.1(A)	Long-term Scheduler.	2-7
2.8.1(B)	Short-term Scheduler	2-7
2.8.1(C)	Medium-term Scheduler	2-7
2.8.1(D)	Comparison of Three Schedulers	2-7
✓	Syllabus Topic : Interprocess Communication.	2-8
2.9	Interprocess Communication	2-8
2.9.1	Message Passing...	2-8
2.9.2	Shared Memory	2-9
✓	Syllabus Topic : Multithreading	2-9
2.10	Multithreading	2-9
2.11	Types of Threads...	2-9
✓	Syllabus Topic : Process - Multithreading Models.	2-9
2.12	Multithreading Models	2-10
z	Syllabus Topic: Thread Libraries....	2-11
2.13	Thread Libraries.	2-11
2.13.1	POSIX Pthreads.	2-12
2.13.2	Win32 Threads	2-12
2.13.3	Java Threads	2-12
z	Syllabus Topic : Threading Issues.....	2-12
2.14	Threading Issues	2-12

2.14.1	The torxu and exec() System Calls	3
2.14.2	Cancellation	3
2.14.3	Signal Handling	3
2.14.4	Thread Pools...	3
2.14.5	Thread-Specific Data	3
2.14.6	Scheduler Activations	3
✓	Syllabus Topic : Process Scheduling• Basic Concepts	3
2.15	Process Scheduling	9
2.15.1	Scheduling Decisions	9
2.15.2	Types of Scheduling	9
✓	Syllabus Topic : Scheduling Criteria	9
2.16	Scheduling	9
2.16.1	Scheduling Criteria	9
✓	Syllabus Topic : Scheduling Algorithms	9
2.16.2	Scheduling Algorithms	9
2.16.2(A)	First In First Out (FIFO)	9
2.16.2(B)	Shortest Job First (SJF)	9
2.16.2(C)	Priority Scheduling	9
2.16.2(D)	Round Robin Scheduling..	9
2.16.2(E)	Multilevel Queue Scheduling	9
2.16.2(F)	Multilevel Feedback-Queue Scheduling	9
2.17	Examples on Uniprocessor Scheduling Algorithms...	9
z	Syllabus Topic : Thread Scheduling	9
2.18.1	Contention Scope	9
2.18.2	Pthread Scheduling	9
z	Syllabus Topic : Multiple Processor Set	9
2.19	Multiple-Processor Scheduling	9
2.19.1	Approaches to Multiple-Processor Sch	9
2.19.2	Processor Affinity	9
2.19.3	Load Balancing	9
2.19.4	Multicore Processors	9
2.19.5	Virtualization and Scheduling	9
2.19.6	Other Multiprocessor Scheduling Applications	9
2.13.6(A)	Load Sharing	9
2.13.6(B)	Sharing Scheduling	9
2.13.6(C)	Dedicated Processor Assignment-	9
2.13.6(D)	Dynamic Scheduling	9
2.20	Exam Pack (University and Previous Questions)	9

Module IP

Chapter 3: Process Coordination 3-1 to 3-27

Syllabns: Synchronization : The critical Section Problem, Peterson's Solution* synchronization Hardware and semaphores. Classic problems of synchronization, monitors, Atomic transactions; Deadlocks System Model, Deadlock Characterization, Methods for Handling Deadlocks, Deadlock Prevention, Deadlock Avoidance , Deadlock Detection, Recovery from Deadlock.

/	Syllabus Topic : Synchronization	3-1
3.1	Background (May 16).....	3-1
3.2	Interprocess Communication (May 16)	3-2
3.3	Race Condition.....*	3-3
z	Syllabus Topic : The Critical Section Problem	3-4
3.4	The Critical Section Problem (Dec. 14, June 15)	3-4
3.5	Mutual Exclusion (June 15, Nov. 15, May 16)	3-5
z	Syllabus Topic : Patersons Solution.....	3-5
3.6	Peterson's Solution.....*.....*	3-5
z	Syllabus Topic ; Synchronization Hardware	3-6
3.7	Synchronization Hardware	3-6
J	Syllabus Topic : Semaphores	3-7
3.8	Semaphores	3-7
z	Syllabus Topic : Classic Problems of Synchronization.....	3-B
3.9	Classic Problems of Synchronization	3-0
3.9.1	Producer Consumer Problem (Dec, 16).....	3-8
3.9.2	Producer Consumer Problem * Using Semaphore (June 15). *	3-8
3.9.3	Readers/Wrtters Problem.....*	3-9
3.9.4	Dining Philosopher Problem (Dec. 14, Nov. 15)	3-10
✓	Syllabus Topic : Monitors	3*11
3.10	Monitors.....*.....*	3-11
z	Syllabus Topic ; Atomic Transactions	3*12
3.11	Atomic Transactions	3-12
3.11*1	System Model of Atomic Transactions	3-12
3.11.2	Log-Based Recovery	3-12
3-11.3	Checkpoints.....	3-13
3.11-4	Concurrent Atomic Transactions	3*13
3.11.4(A)	Serializability.....	3*13
3.11.4(B)	Locking Protocol.....	3*13
3.11.4(C)	Timestamp-Based Protocols.	3*14
z	Syllabus Topic : Deadlocks	3*14
3.12	Deadlocks (Jun* 15, Not* 15* May 16, □. 16)	3-14
z	Syllabus Topic ; System Model	3*14

3.12.1	System Model of Deadlocks	3-14
z	Syllabus Topic : Deadlock Characterization.....*	3-15
3.13	Deadlock Characterzation	3-15
3.13.1	Conditions (June 15, Nov. 15, May 16, Dk. 16).....	3-15
3.13 .2	Resource Allocation Graphs (Dec. 16).....	3-15
z	Syllabus Topic : Deadlock Prevention	3-16
3.14	Deadlock Prevention (June 15, Nov. 15}	3-16
z	Syllabus Topic : Deadlock Avoidance	3-17
3.15	Deadlock Avoidance (June 15, Nov. 15}	3-17
3 15.1	Deadlock Avoidance Algorithms (May 16, Dec. 16)	3-18
3.15.1(A)	Resource-Allocation Graph Algorithm.....	3-19
3.151(B)	Banker's Algorithm(Dec* 16).....	3-19
3.15.1(C)	Resource- Request Algorithm.....	3-19
3.15.1(D)	Safety Algorithm.....*	3-20
3.16	Solved Problems	3-21
z	Syllabus Topic : Deadlock Detection	3-23
3.17	Deadlock Detection (June 15}	3-23
J	Syllabus Topic : Recovery from Deadlock	3-25
3.1B	Deadlock Recovery.....	3-25
3.13.1	Process Termination (Kill a process)-.....*	3-25
3.16.2	Resource Preemption.	3-25
3.19	Exam Pack (University and Review Questions)	3-26

Module IV

Chapter 4 : Memory Management 4-1 to 4-35

Syllabus : Memory Management strategies : Background, Swapping, Contiguous Memory Allocation, Paging , Structure of the Page Table. Segmentation; Virtual Memory Management: Demand Paging. Copy-on- Write, Page Replacement. Allocation of Frames, Thrashing, Memory-Mapped Files, Allocating Kernel Memory , Other Considerations .

z	Syllabus Topic : Memory Management Strategies	4 1
4.1	Memory Management Strategies	4-1
4.1.1	Monoprogramming	4-1
4 1 2	Multiprogramming	4-2
4.1.3	Dynamic Loading	4-2
4.1.4	Overlays	4-2
4.1.5	Relocation	4-3
4.1.6	Logical and Physical Address Space	4-3
✓	Syllabus Topic : Swapping	4-3
4.1.7	Swapping	4-3
✓	Syllabus Topic : Contiguous Memory Allocation.....	4-4
4.2	Contiguous Memory Allocation	4-4

4.2.1	Multiprogramming with Fixed and Variable Partitions.....	4-4
4.2.2	Dynamic Partition Technique	4-5
4.2.3	Compaction.....◆	4-6
4.2.4	Memory Allocation Strategies (May 2016) «.....	4-6
✓	Syllabus Topic : Paging.....	4-7
4.3	Paging (June 15).....j.....	4-7
4.3.1	Basic Operation.....	4-6
4.3.2	Memory Protection and Sharing.....	4-9
4.3.3	* Translation Lookaside Buffer.....1	4-9
4.3.4	Effect of Page Size on Performance (June 15, Nov. 15, Dec. 16)	4-9
4.3.5	Hardware Support for Paging (Dec. 16)	4-9
✓	Syllabus Topic : Structure of the Page Table	4-10
4.4	S. njaure of Page Tables (Dec. 14).....	4-10
4.4.1	Hierarchical Paging	4-10
4.4.2	Hashed Page Table	4-10
4.4.3	Inverted Page Table	11
✓	Syllabus Topic: Segmentation.....	4-12
4.5	Segmentation.....	4-12
4.5.J	Difference b/w P and np a g: seg ~ (Dec. 16)	4-12
4.6	Segmentation p aging.....	4-13
✓	Sy n » Topic : Virtual Memory	4-14
4.7	Virtual Memory Management a ⁹* S	4-15
4.7.1	Virtual Memory (May 16)	MS
✓	S. m-bu. Topic t Demand Paging.....	15
4.7.2	Demand Paging (Dec., 4, Dec. 16)	4-15
4.7.3 S. m-bu. Un and Instruction Restarts;	4-15
✓ S. m-bu. * > *	16
4.0	Ca Py-OH-Write.....	16
✓	Syllabus Topic : Page Replacement.....	4-S
4.9	Page Replace s.....	4-17
4.9.1	FIFO Algorithm May /.....	4-17
4.9 J?	Optimal Page Replacement.....	4-17
4.9.3	LRU Algorithm.....	4-17
4.9.4	• U-Appre, Page' aZ;	4*10
4.9.4(A)	4-1Q
4.9.4(B)	S. m-bu. Chano Algorithm	4-1g
4.9.4(C)	Enhanced SKond -Chance A	4-1g
<9.4(O)	Clock Page flRep lac t At gorithm	4-1g

4.9.5	Counting-Based	Page Replacement	Table of
4.9.5(A)	Not Frequently Used or Little Used	Least Frequently Used Algorithm (LFU or LRU)	Placement
4.9.5(B)	Most Frequently Used Algorithm (MFU).....	Page Replacement	Most Frequently Used Algorithm (MFU or LRU)
4,10	Examples on Page Replacement	Page Replacement	5.1.1
✓	Syllabus Topic : Allocation	Allocation	5.1.2
4.11	Allocation of Frames	Allocation of Frames	5.1.3
✓	Syllabus Topic : Thrashing	Thrashing	5.1.4
<12	■Thrashing(June 15)	Thrashing	5.1.5
4.13	Locality (Working Set Model)	Locality	5.1.6
✓	Syebus Topic =Memory-Mapped Files	Memory-Mapped Files	5.1.8
<14	Memory-Mapped Files	Memory-Mapped Files	5.15(A)
<15	Memory-Mapped I/O	Memory-Mapped I/O	5.1.8(B)
✓	Syllabus Topic : Allocating Kernel Memory	Allocating Kernel Memory	5.1.9
4.16	Allocating Kernel Memory	Allocating Kernel Memory	5.1.10
4.16J	Buddy System.....	Buddy System	5.1.11
<16.1(A)	Operation of Buddy Algorithm	Operation of Buddy Algorithm	5.1.12
<16.2	Slab Allocation	Slab Allocation	5.1.13
✓	Syllabus Topic : Other Considerations	Other Considerations	5.1.13(A)
4.17	Other Considerations	Other Considerations	5.1.13(B)
4.17.1	Prepaging	Prepaging	5.1.14
<17.2	Page Size	Page Size	5.1.14(A)
4.17.3	TLB Reach	TLB Reach	5.1.14(B)
<17.4	Indirect Page Tables.	Indirect Page Tables	5.1.14(C)
4.17.5	Program Structure	Program Structure	5.2
4.18	Program Structure (University and Review Questions)	Program Structure (University and Review Questions)	5.2.1
LMModule V			52.2
5.2.2(A)			52.2
Charging	Storage Management		5.2.2(B)
File System	File System		5.2.2
DetectoX	File System	File Concept.	2.3
Sharing and Protection;	File System	File-System M	2.3(A)
Implementation	File System	File-System M	2.3(B)
hpmemem	File System: Directory	File-System Structure	2.4
File * SpaJ?*	File System: Directory	Implementation.	2.4(A)
ftCC0	File System: NFS	File System: NFS	2.4(B)
Secondary c.	File System: NFS	File System: NFS	2.4(C)
structure. Disk "T	File System: Structure	File System: Structure	2.4(D)
Disk M>.	File System: Structure	File System: Structure	2.4(E)
Implementar	File System: RAID	File System: RAID	2.4(F)
M, 4genCn	File System: Tertiary Storage	File System: Tertiary Storage	2.4(G)
I/O Svcs	File System: StricW®	File System: StricW®	2.4(H)
Kert >elvn V	File System: StricW®	File System: StricW®	2.4(I)
Operations, STREAMS, Performance.	File System: StricW®	File System: StricW®	2.4(J)

✓	Syllabus Topic : Aile System - File Concept.	5-1	5.2.5	Free Space Management	5-25
5.1	File System	5-1	5.2.5(A)	Bit Map or Bit Vector.....	5-25
5.1.1	Pile Concept	5-1	5.2.5(B)	Linked Ust of Disk Blocks	5-26
5.1.2	File Attributes	5-1	5.2.5(C)	Grouping	5-26
5.1.3	Fife Operations	5-2	5.2.5(D)	Counting	5-26
5.1.4	Fite Naming	5-3	✓	Syllabus Topic : Efficiency and Performance	5-26
5.1.5	File Types	5-4	5.2.6	Efficiency and Performance.....	5-26
5.1.6	Fite Structure	5-5	5.2.6(A)	Efficiency	5-26
✓	Syllabus Topic : Access Methods	5-6	5.2.6(B)	Performance	5-27
5.1.7	Access Methods (June 15)	5-6	✓	Syllabus Topic : Recovery	5-28
✓	Syllabus Topic : Directory Structures	5-7	5.2.7	Recovery	5-28
5.1.6	Directory Structures	5-7	5.2.7(A)	Consistency Checking	5-29
5.1.8(A)	Single-Level Directory Systems	5-7	5.2.7(B)	Backup and Restore	5-29
5.1.6(B)	Two-level Directory Systems....	5-8	5.2.7(C)	Log-Structured File Systems	5-29
5.1.8(C)	Hierarchical Directory Systems	5-8	✓	Syllabus Topic : NFS	5-30
5.1.9	Path Names	5-9	5.2.8	NFS	5-30
5.1.10	Directory Operations	5-9	5.2.8(A)	NFS Architecture	5-30
✓	Syllabus Topic : File System Mounting	5-10	5.2.0(B)	NFS Protocols.....	5-30
5.1.11	File System Mounting	5-10	5.2.6(C)	NFS Implementation.....	5-31
5.1.12	Working of Files	5-11	✓	Syllabus Topic : Secondary Storage Structure	5-31
✓	Syllabus Topic : File Sharing	5-12	5.3	Secondary Storage Structure.....	5-31
5.1.13	File Sharing	5-12	✓	Syllabus Topic : Overview of Mass-Storage	5-31
5.1.13(A)	Multiple Users	5-13	Structure	5-31	
5.1.13(B)	Remote File System	5-14	5.3.1	Overview of Mass-Storage Structure	5-31
5.1.13(C) Consistency Semantics	5-14	5.3.1(A)	Magnetic Disks	5-31	
✓	Syllabus Topic ; Protection	5-14	5.3.1(B)	Magnetic Tapes	5-32
5.1.14	Protection	5-14	y'	Syllabus Topic ; Disk Structure	5-32
5.1.14(A)	Type of Accesses	5-15	5.3.2	Disk Structure	5-32
5.1.14(B)	Protection Domains	5-15	✓	Syllabus Topic : Disk Attachment	5-33
5.1.14(C)	Access Control	5-16	5.3.3	Disk Attachment	5-33
✓	Syllabus Topic : File System implementation	5-16	5.3.3(A)	Host- Attached Storage	5-33
52	File System implementation	5-17	5.3.3(B)	Network- Attached Storage (NAS)	5-33
52.1	Fite System Structure	5-17	5.3.3(C)	Storage-Area Network (SAN)	5-33
✓	Syllabus Topic : Implementing File System	5-17	✓	Syllabus Topic : Disk Scheduling	5-34
5.2.2	Implementing Fite System...	5-18	5.3.4	Disk Scheduling (Nov* 15)	5-34
5.2.2(A)	File System Layout..	5-18	5.34(A)	FCFS Scheduling Algorithm (Dec. 16)	5-34
5.2.2(B)	Virtual File System	5-18	5.3.4(B)	Shortest-Seek-Time-First (SSTF) Scheduling	5-34
✓	Syllabus Topic : Directory Implementation	5-19	5.3.4(C)	Algorithm (Dec. 16)	5-34
52.3	Directory Implementation	5-19	5.3.4(D)	SCAN Scheduling Algorithm (Dec. 16)	5-36
5.2.3(A)	Linear List	5-21	5.3.4(E)	C-SCAN Scheduling Algorithm (Dec. 16)	5-36
52.3(B)	Hash Table	5-21	5.3.5	LOOK Scheduling Algorithm (Dec. 16)	5-37
✓	Syllabus Topic : Allocation Methods	5-21	5.3.5	Examples on Disk Scheduling Algorithms	5-42
52.4	Allocation Methods (Dec. 14, May 1®)	5-21	✓	Syllabus Topic ! Disk Management	5-42
5.2.4(A)	Contiguous Allocation.....	5-22	5.3.6	Disk Management.....	5-42
52.4(0)	Linked List Allocation	5-23	5.36(A)	Disk Formatting	5-43
52.4(C)	Linked Ust Allocation using a Table in Memory	5-23	5.36(B)	Boot Block	5-43
52.4(0)	Indexed Allocation	5-23	5.36(C)	Sad Blocks	5-43
5.2.4(E)	I-nodes,..	5-25	✓	Syllabus Topic ; RAID Structure	5-43
✓	Syllabus Topic £Free Space Management.....	5-25	5.3*7	RAID Structure	5-43

5.37(A)	RAID Levels.....	5-43
✓	Syllabus Topic : Stable-Storage Implementation	5-46
5.3.6	Stable-Storage X'Xstorag® Structure.....	5-46
✓	Syllabus Topic : Tertiary-Storage Structure	5-46
5.3.9	Te «-S» Devces.....	2
5.3.9(A)	Te «-S» Devces.....	5-48
5.3.9(B)	Osra iing-sen, Suppo pedomance TM.....	5-49
5.3.9(C)	ace Management.....	5-49
✓	Syllabus Topic : Swap-w	5-49
5.3.10	swa p-SpaceMaW Im9nl.....	5-49
5.3.10(A)	Swat>Spa [»]	5-50
5.310(B)	SwafhSpaM.....	5-50
	Sy1M>b. Topic : I/O systems	5-50
5.4	VO systems.....	5-50
5.4.1	Overview.....	5-50
5.4.1(A)	I/O Devices,™.....	5-50
5.4.1(B)	Differences between VO Devices.....	5-51
✓	Syllabus Topic : Overview VO Hardware	5-51
5.4.2	VOHardware.....	5-51
5.4.2(A)	Device Controllers.....	5-51
5.4.2(B)	Polling.....	5-51
5.4.2(C)	interrupt Handier.....	5-52
5.4.2(D)	Interrupt Service Routine (ISR).....	5-53
5.4.2(E)	Direct Memory Access (DMA).....	5-53
✓	Syllabus Topic : Application I/O Interlace	5-54
5.4.3	Application 1.0 Interlace.....	5-54
5.4.3(A)	Block or Character Devices.....	5-55
5.4.3(B)	Network Devices.....	5-55
5.4.3(C)	Clocks and Timers.....	5-55
5.4.3(D)	Blocking and Non-blocking VO.....	5-55
✓	Syllabus Topic : Kernel I/O Subsystem	5-56
5.4.4	Kamel VO Subsystem.....	5-56
5.4.4(A)	VOScheduling.....	5-56
5.4.4(B)	Buffering,™.....	5-56
5.4.4(C)	Caching.....	5-57
5.4.4(D)	Spooling and Device Reservation.....	5-57
5.4.4(E)	Error Handling.....	5-57
5.4.5	VOProtection.....	5-58
5.4.6	Kamel Data Structures.....	5-58
✓	Syllabus Topic : Transforming I/O Requests to Hardware Operations	5-58
5.4.5	Transforming I/O Requests to Hardware Operations.....	5-58
✓	Syllabus Topic: STREAMS...	5-59
5.4.6	STREAMS.....	5-59
✓	Syllabus Topic : Performance	5-59
5.4.7	Performance.....	5-60
5.5	Exam Pack (Untamty and Review Questions).....	5-60
		5-60

Chapter 6 : Distributed Systems

✓	Syllabus Topic : Network Communication	5-8.1
	System : Network Communication	5-8.2
	Naming and	3
	Protocols	3
	Distributed Systems	3
	file access. Stateful Versus Stateless, X	3
	Locata	3
	Distributed Fanrl Upaflocl	3
	Concurren	3
	j Trope: Distributed Opeptfngs	3
✓	Distributed Systems.....	1 A3
6.1	Definition	9.4
6.1.1	Definition	9.4
6.1.2	Motivation	9.5
6.1.2(A)	Resource Sharing.....	3
6.1.2(B)	Computation Speedup	10
6.1.2(C)	Reliability	11
6.1.2(D)	Communication	11
6.2	Types of Distributed Operating Systems.....	12
✓	Syllabus Topic ; Network Based OS-TM,™.....	1 Mi
6.2.1	Network Operating System (NOS).....	12.1(A) Ce
6.2.1(A)	Remote Login	12.1(B) Fuf
6.2.1(B)	Remote File Transfer	12.1(C) To*
6.2.2	Distributed Operating Systems (DOS).....	
6.2.2(A)	Data Migration	
6.2.2(B)	Computation Migration	
6.2.2(C)	Process Migration	
✓	Syllabus Topic : Network Structure	
6.3	Network Structure	
6.3.1	Local-Area Networks (LANs)	
6.3.2	Wide-Area Networks (WANs)	
✓	Syllabus Topic : Network Topology	
6.4	Network Topology.....	
✓	Syllabus Topic : Communication Structure	
6.5	Communication Structure	
6.5.1	Naming and Name Resolution.....	
65,2	touting Strategies	
65,3	Jacket Strategies.....	
6.5.4	Connection Strategies	
6.5.5	Contention.....	
✓	Syllabus Topic : Communication Protocols	
6.6	Convocation Protocols	
✓	Syllabus Topic : Distributed HI@SV5	
6.7	DhWb «W FH Sy8t8ms.....	
	Syllabus Topic : Naming and Transparency	

6.8	Naming and Transparency	6-10	v	Syllabus Topic : Concurrency Control	6-16
6.8.1	Naming Structures.....	6-1 0	6.13	Concurrency Control	6-16
6.8.2	Naming Schemes	6-1 1	6.13.1	Locking Protocols	6-16
6.8.3	Implementation Techniques	8-1 1	6.13.1(A)	Nonreplicated Scheme	b-16
z	Syllabus Topic : Remote File Access	6-12	6.13.1(8)	Single-Coor dinato r Approach	6-16
6.9	Remote File Access	6-12	6.13.1(C)	Majority Protocol	6-17
6.9.1	Basic Caching Scheme	6-12	6.13.1(D)	Biased Protocol	6-17
6.9.2	Cache Location.....	6-12	6.13.1(E)	Primary Copy	6-17
6.9.3	Cache-Update Policy	6-12	6.13.2	Timestamping	6-17
6.9.4	Consistency...,	6-13	6.13.2(A)	Generation of Unique Timestamps	6-17
6.9.5	A Comparison of Caching and Remote Service.	6-13	6.13.2(B)	Timestamp-Ordering Scheme..	6-18
✓	Syllabus Topic : Stateful Versus		J	Syllabus Topic : Deadlock Handling	6-18
	Stateless Service.....	6-14	6.14	Deadlock Handling	6-18
6.10	Stateful Versus Stateless Service	6-1 4	6.14.1	Deadlock Prevention and Avoidance.	6-18
✓	Syllabus Topic : File Replication	6-14	6.14.2	Deadlock Detection....	6-19
6.11	Ale Replication..	6-14	6.14.2(A)	Centralized Approach	6-19
✓	Syllabus Topic : Distributed Synchronization	6-15	6.14.2(B)	Fully Distributed Approach....	6-19
6.12	Distributed Synchronization.	6-1 5	6.15	Exam Pack (Review Questions)	6-20
6.12.1	Mutual Exclusion	6-15	•	Lab Manual	L T to L-26
6.12.1(A)	Centralized Approach	6-15			
6.12.1(B)	Fully Distributed Approach	6-15			
6.12.1(C)	Token-Passing Approach.,	6-16			



Overview of Operating System

~ ⑧y1,abua J?P'S..

Introduction : Operating System Structure

storage management, Protection and security

Structure: Operating system services and interface

Operating System Design and ZtemeZ ion Os' X,

generation. System boot .

ZT" 1 Memory mana geme "k

P Ca purpose

Systems; System

X T uis types. Sy,em

rams,

structure. Virtual machines, OS debugging and

Syllabus Topic : Introduction

1.1 Introduction to Operating System

What is OS?

→ (Dec. 14, June 15, Nov. 15)

Q- What is operating system?

MU - Dec. 2014, June 2015, Nov. 2015, 2 Marks

An operating system is system software which manages, operates and communicates with the computer hardware and software. To complete the execution user program need many resources. The main job of the **operating system is to provide resources and services to the user program**. So without operating system, a computer would be useless.

Definition of operating system

- An operating system acts as an interface between the user and hardware of the computer and also controls execution of application programs.

- Operating system is also called as resource manager.

Need of Operating System

- Basically operating systems identifying performing tasks, for example keyboard and the input devices such as devices such as monitor, printer and keyboard and output of files and disk. Or, printer and keyboard and disk on the disk. **storage devices**, peripheral devices, Printers, scanners, audio mixers. The **called CPU** system is a processing unit **program should get the processing**

unit to complete the execution and computer system should offer the service to allocate processing unit to users program.

- The operating system allocates memory to user program as per need.
- In the same way, user programs interact with the other user programs through devices like key board or a mouse or even a joy stick.

Syllabus Topic : Operating System Structure

1.2 Operating System Structure

- As single user cannot keep CPU or I/O devices busy all the time, operating system supports multiprogramming.
- In multiprogramming, operating system organizes the jobs in such a way that the CPU always get one job to execute, Hence, Multiprogramming improves CPU utilization.
- Multiple jobs remain in memory so that CPU can take any one for processing.
- During processing of job if I/O is required then CPU does not remain idle it takes other job for processing till I/O completes.
- Multiple users are permitted by time shared operating system to share the machine at the same time.
- Time shared system is capable of completing each action or command in a very less time. As a result, each user need can be fulfilled in very less time.
- A particular user program executes in allocated time by system and after expiration of slot, system gets allocated to other user program.

1ST Operating System (MU - Sem 4) T1

- So the allocation of computer resources in time slots to several programs simultaneously is in time sharing manner. Consider the example of multiple use in same server.
- The resources in server machine are allocated to each user on time sharing basis giving the feeling to each user that only he or she is using server machine exclusively.

Designs of the operating system

- Following are some of the designs of the operating system that is tried practically.

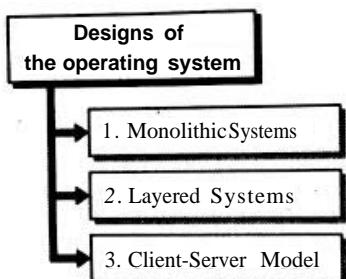


Fig. C.1.1 : Designs of operating system

-F 1.2 J Monolithic Systems

Q. Explain monolithic system.

- This type of operating system can be treated as having no structure.
- The operating system is constructed as a set of procedures and each procedure can call any other ones whenever needed.
- Each procedure in the system has a precise interface in terms of parameters and results.
- Every procedure can call any other one, if the called one offers some useful computation that the calling procedure needs, caning
- In this technique, compilation of all the procedures or files containing the procedures is <2

- Every procedure is able to be seen as a procedure so there is no function "very other
- A little structure is possible to implement in monolithic systems,
- The system calls offered by the operating system are requested by putting the parameters in the stack and then place (e.g., on the stack) and then instruction. cutting a trap
- Due to this instruction machine gets SV mode to kernel mode and transfers control to the operating system.
- The operating system then fetches the find out which system call is to be execut

- This organic operating system suggests a basic structure for the system:
- A main program that calls several procedures that are demanded by the system.
 - A Xwn of service procedures that perform system calls.
 - A collection of utility procedures that assist the service procedures.
 - In this representation, there is a single service procedure for each system call which takes care of it.
 - The utility procedures perform things that are done by several service procedures, such as fetching from user programs.
 - This separation of the procedures into three layers is shown in Fig. 1.2.1.

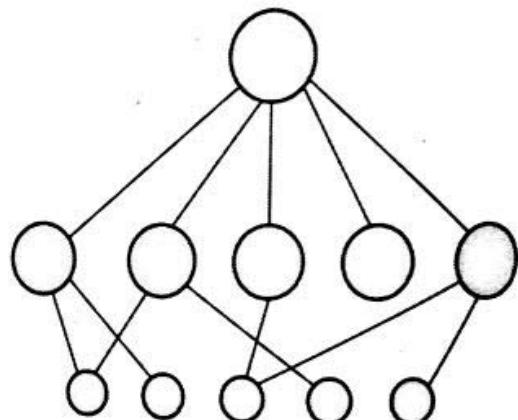


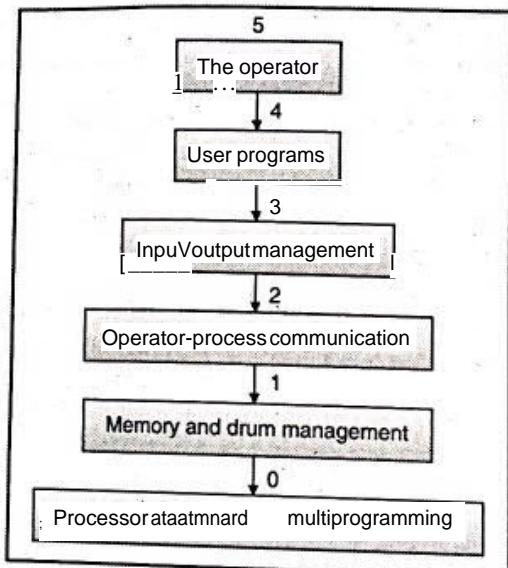
Fig. 1.2.1: A simple structuring model for a monolithic system

→ 1.2.2 Layered Systems

Q. Explain layered system in detail.

- Another approach to the structure of the system is to build it upon the layers. One constructed
- The first system built in this way was the THE system developed by the N. K. Rijf, Technische Hog school Eindhoven, students, and W. Dijkstra (1968)
- As shown in Fig. 1.2.2, the system had total 6 layers.
- Layer 0: This layer was the "processor, controller, SpOBSiWe" to handle allocation of memory when interrupts occur.
- Layer 1: This layer is as responsible to offer memory management of the CPU.
- Layer 2: This layer was responsible to allocate space for processes.
- Layer 3: This layer was responsible to handle memory management.
- Layer 4: This layer was responsible to handle memory management.
- Layer 5: This layer was responsible to handle memory management.

- If (here is no space in main memory, then it also allocates space on a 512K word drum used for holding parts of processes(pages)).
 - The layer 1 software was concerned of making needed pages were brought into memory whenever they were required,



Hg* 1.2.2 : Structure of the THE operating system

☞ Layer 2

- " " ka,ion betw" n each PTOe" and,he
oXr con'X e, Ch Pr " eSS effcCliVCly had ,S Own

☞ Layer 3

- P* management of I/O devices and buffering the streams and from them was handled «* b y
 - **On top of layer 3 each process could be able to deal with many peculiarities.**

Layer 4

- The user programs were resided in this layer

mer_{XeT₂}

Layers

astern XulTT o'd o f the WayS 10 make
Fig- 1.23, the onenf approach shown in

The bottom layer is hardware and the highest layer is user interface.

- An operating system layer is an implementation of an abstract object made up of data and the operations that can manipulate those data.
 - A characteristic operating system layer consists of data structures and a set of routines that can be invoked by higher-level layers.
 - The same layer can too invoke operations on lower-level layers.

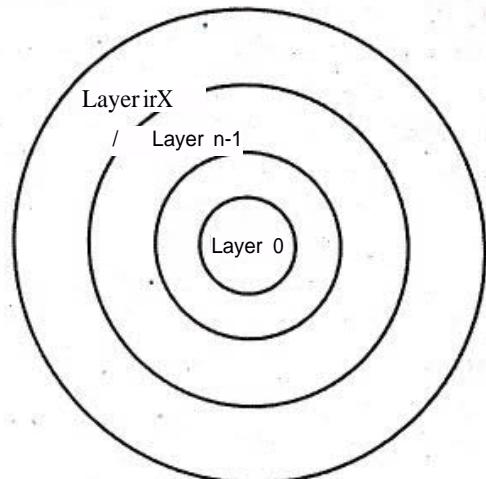


Fig. 1.23 : A layered operating system

→ Advantages

The advantages of layered approach are that

- It keeps much better control over the computer and ~~over the applications that utilizes that computer.~~
 - Implementors have more liberty in changing the internal workings of the system and in writing modular operating systems.
 - **Ease of construction and debugging.**

☞ **Limitations of layered approach**

- defin vari ~~tr_o~~Ubk W,th the layered a PProach is to
lower-level layers, cautious planning is necessary **only**

■ *1.2.3 Client-Server Model

Q. Explain client-server model in detail.

- As a large part of the traditional operating system code of ~~simX~~ ~~high-level CMS~~, ~~U1e VM/370 g0tal0t~~
 - ~~However VM/370~~ ~~is a~~ ~~task~~ ~~is~~ ~~complex~~ ~~task~~
 - A trend in modern operating systems is to take the idea of moving code up into higher layers even further and



- remove as much as possible from kernel mode, leaving a minimal microkernel.
- The usual approach is to implement most of the operating system in user processes.
- To request a service, such as reading a block of a file, a user process (client process) sends the request to a server process, which then does the work and sends back the answer.
- In this model, shown in Fig 1.2.4, all the kernel handles the communication between clients and servers.
- By splitting the operating system up into parts, each of which only handles one facet of the system, such as file service, process service, terminal service, or memory service, each part becomes small and manageable.
- As all servers run as user-mode processes, and not in kernel mode, they do not have direct access to the hardware.
- As a consequence, if a bug in the file server is triggered, the file service may crash, but this will not usually bring the whole machine down.
- Another advantage of the client-server model is its adaptability to use in distributed systems.
- **The communication of client with a server** takes place by sending messages.
- When client sends a message to server, it is not necessary for client to know whether the message is processed in its own local machine or sent to a server on a remote machine in network.
- Either server is present on client machine or remote machine in network from client's point of view. It appears to be a request was sent to the server and a response came back.

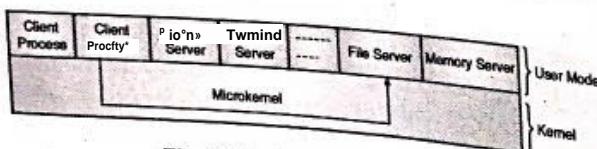


Fig. 1.2.4. Client-Server Model

1-2.4 Monolithic Kernel vs. Microkernel

Q. Differentiate between monolithic and microkernel (Jnu 15, Nov. 15, Dec. 16)

Q.7 What is Kernel? Describe briefly the approaches of designing Kernel. (MU - Dec. 2015, 5 Marks)

- The majority of CPUs support two modes: kernel mode and user mode, i.e. kernel mode and user mode.
- In kernel mode, all instructions are executed, and the entire memory is accessible throughout the system.

- On the contrary to this, in user mode, register access is restricted.
- The CPU gets switched to kernel mode while operating system code.
- The only means to switch from user mode to kernel mode is through system calls as implemented by the operating system.
- The operating system can be put in full system calls are the only basic services an OS provides.
- The operating system can also be put in full system calls are the only basic services an OS provides.
- If virtually entire operating system code is executed in kernel mode, then it is a monolithic program in a single address space.
- The disadvantage of this approach is that it is difficult to change or adapt operating system components without doing a total shutdown and perhaps even recompilation and reinstallation.
- Monolithic operating systems have drawbacks from the viewpoint of openness, software reliability, or maintainability.
- If the operating system is organized into two parts, it can provide more flexibility.
- In the first part, the hardware is kept which can uniformly well be used in user mode.
- For example, the memory management module keeps allocated and free space to the processes.
- It is required to execute in kernel mode at the time of the MMU reset. of the Minor memory Management.
- execute in user mode. It is only code in the operating system.
- Actually a for context switching, it requires only containing the MMU. Setting device, MMU, and capturing.
- Microcode is used to pass system calls to the proper user level operating system.
- It is used to switch their results. system with the organization of operating system.

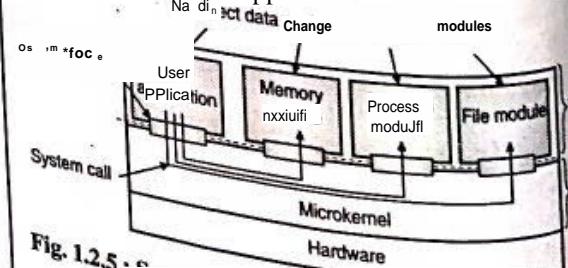


Fig. 1.2.5. Monolithic Kernel

Urgent for operating system as a microkernel

**Difference between monolithic and microkernel**

Q. Differentiate between monolithic and microkernel

→ (June 15, Nov, 15)

MU-June 2015, Nov. 2015. 5 Marks

Sr. No.	Parameters	Microkernel	Monolithic kernel
1.	Definition	In microkernel, set of modules for managing the hardware is kept which can uniformly well be executed in user mode. A small microkernel contains only code that must execute in user mode. <u>End mode is the part of the operating system</u>	If virtually entire operating system code is executed in kernel mode, then it is a monolithic program that runs in a single address space
2.	Address Space	There is separate address space for user services and kernel services.	There is same address space for user services and kernel services
3.			
4.	Execution speed	Microkernel size is smaller than monolithic kernel. Execution speed is slower than monolithic kernel.	It has a larger size compared to microkernel
5.	Flexible	It is more flexible. Any module can be replaced.	Execution speed is faster than microkernel.
6.	Crash	Crash of one service does not affect working of other part of microkernel.	It is not flexible as compared to microkernel. Its component cannot be changed, without doing a total shutdown and perhaps even a full recompilation and reinstallation.
7.	Communication	In microkernel, communication is through messagepassing.	After crash of service complete operating system fails.
8.	Drawbacks	Performance degradation due to context switching and communication overhead between different modules.	Kernel calls the function directly.
			Monolithic operating systems have drawbacks from the viewpoint of openness, software engineering, reliability, or maintainability. It has lack of flexibility and having huge size.

Syllabus Topic j Operating System Operations

1.3 Operating System Operations

Q. Byain operating system operations.

- Modern operating systems are interrupt driven. Operating system sits idle and waits in case there are no devices to offer service. Given no USC to whom response is to be given.
- In this case operating system waits for something to happen. When interrupt or trap occurs then events are triggered. If detected by CPU is called trap.
- Trap or exactly same point of time of program execution. Software interrupts are trap or exceptions. Example its division by zero.

- After receiving an interrupt, operating system can out some housekeeping so that it can resume its computation once it is through servicing the interrupt.
- After this a searching is carried out in the interrupt vector or interrupt table.
- This table remains in kernel-memory space and includes address of the code in the device driver that is to service the interrupt.
- The interrupt handler is then executed. When the handler finishes, control of the CPU is returned to the code that was executing before the interrupt occurred.
- Interrupt service routine (ISR) deal with interrupts.
- The design of the operating system should be in such a way that one erroneous or malicious program should not affect other programs.

1.3.1 Dual Mode Operation

- There are two modes of operating system.
- The majority of CPUs support at least two modes of operation i.e. kernel mode and user mode.

- Operating System (1)

allowed to be t of all registers

 - „K.™, - M - - - » executed, and the entire memory an sc is accessible throughout the execution.
 - On the contrary is to this in **user mode**, memory and to register access is restricted. The CPU gets switched to kernel mode while executing operating system code.
 - The only means to switch from user mode to kernel mode is through system calls as implemented by operating system.
 - Hardware contains mode bit which indicate kernel mode when 0 and user mode when set to 1. When user application is running then system is in user mode.
 - When user application requests the service from operating system then transition from user mode to kernel mode takes place.
 - During booting of the system, hardware is in kernel mode. Once the operating system is loaded and starts executing user application system goes in user mode.
 - Due to dual mode operation, protection from misbehaving users is achieved. This protection can be provided by allowing the execution of privileged instructions in kernel mode only.
 - Any attempt to execute these privileged instructions in user mode causes a trap to the operating system.
 - The instruction to enter user mode, instructions for I/O control, timer management, and interrupt management are some of the examples of privileged instructions.
 - From user mode, control is switched to operating system through interrupt, a trap, or a system call.
 - The interface between OS and user programs is defined by a set of system calls that the operating system provides.
 - System call is the call for the operating system to

- There are two ways for the n_{roc} to switch from the user mode to the kernel mode.
 - One is when a user process explicitly requests¹⁰ kernel mode by issuing a system call.
 - Other is during the transition of user process kernel can take over to system "housekeeping task.

13.2 Timer

- Computer contains ~~sta~~ | clocks containing crys oscillator, a counter, and a holding register.
 - A quartz crystal under tension produces a periodic ~~very high~~ accuracy, usually in the range of X_1 — * megahertz to a few Gigahertz, depending on the crystal chosen.
 - By electronic circuitry, this base signal can be multiplied by a small integer to get frequencies up to several gigahertz and more than this.
 - Computer contains at least one circuit, than produce synchronizing signal and given to many circuits in computer.
 - This synchronizing signal is fed into the counter to make it count down to zero. When the counter gets to zero, it causes a CPU interrupt.
 - In one-shot mode, after the clock is started, holding register value gets copied into the counter and then decrements the counter at each pulse from the crystal. I
 - When the counter value becomes zero, it causes an interrupt and stops until it is explicitly started again by the software.
 - In square-wave mode, after getting to zero and causing the interrupt, the holding register is automatically copied into the counter, and the whole process is repeated again indefinitely.
 - These periodic interrupts are called clock ticks. Every computer has a battery-powered backup clock to save the current time when machine is switched off.
 - Universal coordinated time is used to synchronize the docks of the machines.
 - Operating system ensures to ~~set timer to~~ in terni ~~before~~ hand over the control to user application.

1.4

-- P - e | s Management

- Q. Explain process ~ " operating system function

A. program or n ...
P . A process execution is called
a program re des *** execution context j
require any de, on * disk. On disk it do* "1
resource



- A program gets executed in main memory. So it should be transferred from disk to main memory.
- To complete execution, program needs many resources, ~~and competes for it. These resources are memory, files etc~~
- now it becomes process. From the computation context X oMiew. a process is defined by CPU state, memory contents and execution environment,
- Process manager implements the process management ~~functions~~.
- tomnWrogramming. single CPU is shared among
- k X processes remain busy in completing UO CPU is allocated to only one process at a given point of time.
- Here ~~some~~ policy is required to allocate CPU to J ess, called as CPU scheduling.
- If multiple users ~~are~~ working on the system, operating system ~~switches~~ the CPU from one user process to ~~other~~
- User gets the illusion that only he or she is using the system.
- Process synchronization mechanism is required ~~critical section~~ to ensure that only one process should use
- Process ~~communication and synchronization, deadlock handling, suspension~~ and resumption of processes and ~~creation and deletion~~ of the processes etc are some of the activities performed in process management.
- ~~Program counter~~ contains address of the instruction from which execution resumes after context switch.
- The process management activities involves .
 1. To provide control access to shared resources like file, memory, I/O and CPU.
 2. Control the execution of user applications,
 3. Creation, execution and deletion of user and system processes.
 4. Resume a process execution or cancel it.
 5. Scheduling of a process.
 6. Synchronization, interposes communication and deadlock handling for processes.

Syllabus Topic : Memory Management

1.5 Memory Management

Q. Explain memory management function of operating system.

- Memory is an important resource of the computer system that needs to be managed by the operating

system. To execute the program, user needs to keep the program in main memory.

The main memory is volatile. Therefore a user needs to store his program in some secondary storage which is non volatile.

- Every process needs main memory since a process code, stack, heap (dynamically- allocated structures), and data (variables) must all reside in memory.
- The management of main memory is required to support for multiprogramming. Many executables processes exist in main memory at any given time.
- Different processes in main memory have different address space.
- Memory manager is the module of the operating system ~~responsible for managing memory~~.
- Programs after completion of execution move out of main memory or processes suspended for completion of I/O can be swapped out on ~~secondary storage to the~~ main memory for other processes.
- New processes are required to be loaded in main ~~memory if available main memory is not sufficient to~~ hold all processes swapping between main memory and secondary memory is done.
- Memory managers move processes back and forth between memory and disk during execution.
- So it is required that operating system should have some strategy for the management of memory.
- The memory management activities handled by operating system are :
 1. Allocation of memory to the processes.
 2. Free the memory from process after completion of execution.
 3. Reallocation of memory to a program after used block becomes free.
 4. Keep track of memory usage by the process.

** Syllabus Topic : Storage Management

1.6 Storage Management

- The operating system ~~gives a uniform, logical view of information storage for making the use computer system convenient to users.~~



- The OS abstracts from the physical storage unit, the storage devices to describe a logical file. • Told out by operating system
- The mapping of files is carried out by the operating system onto physical media. These files are accessed by operating system via the storage devices.

1.6.1 File Management

Q. Explain file management in OS.

- Following are the necessities for long term information storage,
 - c It must be possible to store a very large amount of information.
 - o The information should not loss after termination of the process using it.
 - o Several processes must be able to access the information simultaneously.
- In order to fulfill the above requirements, it is necessary to store information on disks and other secondary storage in units called **files**.
- A file is a named collection of related information that is recorded on secondary storage.
- The data cannot be written directly to secondary storage unless they are within a file. Processes can then read information from the file. It also can write new information into the file if needed.
- After process termination, information in file should remain retained and should not vanish.
- A file should only vanish when its holder clearly tells the system to do so.
- File system describes how files are structured, named, accessed, used, and protected.
- The file management activities of the operating system consist of:
 - 1. Creation and deletion of files and directories.
 - 2. Provide access to files and directories.
 - 3. Management of storage space for files.
 - 4. File Security.
 - 5. Providing interfaces for file management.

1.6.2 Mass Storage Management

Q. Write note on mass storage management.

j, volatile. Storage devices and Main memory devices are used in computer system. Storage devices are ssion devices (network cards (disks, tapes), modems), and human-interface keyboard, mouse).

Two main types of storage technology used today are magnetic and optical storage. Primary magnetic

“Alices are: Diskettes, Hard disks (both fixed and removable), High capacity floppy disks, cartridges, and Magnetic tape.

Primary optical storage devices are: Compact Disk, Read Only Memory (CD ROM), Digital Video Disk, Read Only Memory (DVD ROM), CD Recordable (CD R), CD Rewritable (CD RW), Photo CD.

The hard disk is used as main storage device in your computer system. It is magnetic type of storage device.

Within one hard disk unit, many physical disks are present. Each disk is known as a platter.

Several metal platters are present within the hard disk and these are coated with a special magnetic material.

The platters rotate many thousands of times within a second for accessing the data. Magnetic read and write head are present which floats just above the surface of the platter.

Material used for making the platters is aluminum, glass, or ceramic and two read/write heads are present, one for upper and lower surface.

Platters are arranged in stack because of which the stack is referred to as cylinder.

• its cylinder has often 1s referred to as cylinder.

Cylinder is the location of a single track through all platters. Due to this, the read and write mechanism is time required to move to hard disk is improved.

Now days hard disk is partitioned and of hard disk can be used as separate drive. Storage capacity of hard disk can be in gigabytes.

The device management tasks include :

1. To open, close and write device drivers,
2. Communicate, control and monitor the device driver.

1.6.3 Caching

Q.

“Wain how caching is implemented.

The memory system can be viewed as hierarchy of 4 layers. Top layer contains many registers which are inside the CPU.

As material used to build the registers and CPU is same, these registers are faster just like CPU.

These registers can be accessed quickly. For these registers, typical access time is around 1 nanosecond and storage capacity is around 1KB.

The second layer of the hierarchy is cache memory. The frequently used data and instructions are cached in cache memory to improve the performance. This memory is controlled by the hardware.

The division of primary memory is done into cache lines of 64 bytes, having addresses 0 to 63 in cache line 0, 64 to 127 in cache line 1, and so on. The most frequently used cache lines are kept in a high-speed cache sited inside or very close to the CPU.

If cache hardware finds that, needed line by program is in cache then it is cache hit.

If it is not in cache then it is treated as cache miss and needs to pay for more time as request would be satisfied from main memory.

For cache, typical access time is around 2 nanosecond and storage capacity is around 4 MB.

In main memory, a buffer is reserved for disk sectors called as disk cache.

The cache keeps a copy of data in some of the sectors on the disk.

To fulfill an I/O request for a particular sector, first disk cache is checked to see whether the sector is in disk cache.

If it is present, then request is fulfilled via the cache. If not present, then required sector is read into the disk cache from the disk.

The disk cache improves the performance as some request is satisfied from it.

The property of locality of reference is used. When a block of data is present in cache to fulfill a single I/O request, it is expected that same block will be needed in future.

Due to disk cache mechanism, time required to read from and write to hard disk is improved. Now days it is the part of hard disk.

Hard disk can be partitioned, and each partition can be used as separate drive. Storage capacity of hard disk

can be in gigabyte or terabytes. For magnetic disk, typical access time is around 10 nanosecond and storage capacity is around 1 to 4 TB.

The recently read data remain in disk cache. Sometimes adjacent data which is expected to be accessed next time is held by disk cache. Some disk cache offers write caching also.

1.6.4 I/O Systems

- As we know that I/O is one of the main functions of an operating system and is used to control the entire computer's Input/Output devices.
- Basically it should have to issue commands to the devices, catch interrupts, and handle errors.
- It should also provide an interface between the devices and the rest of the system that is simple and easy to use.
- The I/O code represents a significant fraction of the total operating system. Computers operate with many kinds of devices.
- As we know that it include storage devices (disks, tapes), transmission devices (network cards, modems), and human-interface devices (screen, keyboard, mouse).
- Following are the components of I/O subsystem.
 - o Buffering, Caching, and Spooling which are memory management components.
 - o Device driver interface.
 - o Drivers for specific devices.

"s n busTopicTprotection and Security

1.7 Protection and Security

Q. Explain protection and security mechanism of OS.

- Protection mechanisms refer to the particular operating system mechanisms which are used to protect information, files and resources in the computer.
- Policy means whose data should be protected from whom and mechanism is how system put into effect these policies.
- In some systems, a program called reference monitor is used to impose a protection.
- Any attempt to access the resource is verified by reference monitor whether it is legal or not.



The referee makes a decision on the basis of policy

- table.
- Modern protection mechanism developed to improve the reliability of « Plex system in which use of shared resources is involved. There are many reasons to offer the protection to avoid protection of an access.
- It is needed to provide the protection of an access to mischievous, deliberate violation.
- It is a must to make sure that, each program component running in the system uses system of the system per defined policies to ensure the reliability.
- Policies to make use of the resources of computer system are put into effect by mechanism which is offered by protection.
- Some of the policies are included in the design of the system. Other policies are decided by the management of a system.
- Many policies are defined by the individual user of the system to protect their own files and programs.
- A protection system should be flexible enough to offer different and put into effect the different types of policies.
- **Different types of applications have different types of resource use.**
 - The policies designed should allow to offer this resource use need.
 - The resource use of the applications can change over the period of time.
 - Hence, instead of relying totally on operating system, application programmer should use protection mechanism to protect the resources created against misuse.
 - The system is secure if its resources are used and accessed as proposed under all conditions. Total security is impossible to achieve.
 - We must have mechanisms to make security breaches an exceptional occurrence, rather than the rule.
 - The challenge in developing operating systems security is to design security mechanisms that safeguard user execution and their generated data in an environment with complex interactions.
 - While protection of the system related to internal environment, security, on the other hand, consideration of the external environment, which

the system operates.

- Only sufficient protection system is not enough for total security. In particular operation protection needs to be used to protect the system. The computer security encompasses the overall problem of security.
- Security has many aspects. There are more ones.

The nature of these threats,

- 0 The nature of intruders.
- 0 Accidental data loss.
- 0

1. Threats

From the viewpoint of security, computer systems have three general goals. With corresponding threats to them,

These are shown in fig. 1.2

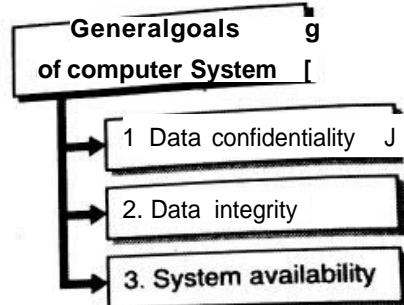


Fig. C1.2 : General goals of computer systems

→ 1. Data confidentiality

- Threat to this goal is publicity of data. All confidential data should remain confidential.
- System should always ensure that, the access rights to data should be given only to authorized person.
- At least, the owner of the data should be able to state who can access what, and the system should enforce these specifications.

→ 2. Data integrity

- Threat to this goal is modification of data. Unauthorized users should not be able to make modification to any data without the owner's permission.
- Data modification includes not only changing the data, but also deleting data and addition of false data as well.
- The system should always guarantee that data stored in it, remains unchanged until the owner decides to change it.

4 3. System availability

- System should be usable at any time. Nobody should disturb the system to make it unusable.
- Denial of service attack is increasingly common to make system unavailable.
- Protecting individuals from misuse of information about them is called privacy.
- This quickly gets into many legal and moral issues. System should also ensure the privacy to individual user.
- The system contains many hardware objects such as CPUs, memory segments, disk drives, printers, magnetic tapes and many software objects such as processes, files, databases, or semaphores.
- Each object can be referenced by its unique name. On every object a finite set of operations can be performed. For example, WAIT and SIGNAL on semaphores and READ and WRITE on files.
- A system should enforce a mechanism, to restrict the processes from accessing the needed objects for which they are not unauthorized.
- The mechanism should also ensure to restrict processes to a subset of the legal operations when that is needed.
- For example, process P has a permission to read file F but not of performing write operation on it.
- A set of object and rights pair is called **domain**. Each pair denotes an object and some subset of the operations that can be performed on it.
- One domain corresponds to one user and specify the permissions to user for certain activities. Consider the following three domains.
- It is possible for the same object to be in multiple domains. [Read, Write, execute] rights are available on each object. At a particular time of execution, each process executes in some protection domain.
- In that domain there is some set of objects it can access, and for each object it has some set of rights shown in square brackets.
- During execution, processes can go from one domain to other domain. The rules for domain switching are very much depends on and varies from system to system
- In Unix every process is defined by user-id and group-id(uid,gid).
- Two processes with the similar (UID, GID) combination will have access to precisely the same set of objects.

Overview of Operating System

- Process often switches from user area to kernel area, that is, from one domain to other.
- The kernel part has access to a different objects ~~user part~~.

Sy,tabua Topic : Distributed Sya~

1.8 Distributed System

- A computer network is defined as a set of communicating devices that are connected together by communication links.
- These devices include computers, printers and other devices capable of sending and/or receiving information from other devices on the network.
- These devices often called as node in the network. So computer network is interconnected set of autonomous computers.
- *A distributed system is defined as set of autonomous computers that appears to its users as a single coherent system.* Users of distributed system feel that, they are working with a single system.
- Distributed system is like multicomputers spread worldwide. Each node in distributed system is having its own CPU, RAM, network board, OS, and disk for paging.

Main characteristics of distributed system

Q. What are the characteristics of distributed system? Explain.

- A distributed system comprises computers with distinct architecture and data representation. These dissimilarities and the ways all these machines communicate are hidden from users.
- The manner in which distributed system is organized internally is also hidden from the users of the distributed system.
- The interaction of users and applications with distributed system is in consistent and identical way, in spite of where and when interaction occurs.
- A distributed system should allow for scaling it.
- Distributed system should support for availability. It should be always available to the users and applications in spite of failures.
- Failure handling should be hidden from users and applications.

JdPF Operating System (MU-Se T IT) user

Comparison of multi-processing, multi-computing and distributed system

Parallel		Multi-processing	Multi-computing	Distributed System
Node configuration	CPU	CPU, RAM, net interface	Computers, input	
Node peripherals	All shared	Shared memory, maybe disk	Full system per node	
Location	Same rack	Same room	Possibly worldwide	
Internode communication	Shared HAM	Dedicated interconnect	Traditional network	
Operating systems	One, shared	Multiple, same	Possibly different	
Filesystems	One, shared	One, shared	Each node has own	
Administration	One organization	One organization	Many organizations	

- Network contains computers with different architectures (heterogeneous).
- **To offer the single system view of distributed system** in this heterogeneous environment, a middleware layer is often placed between higher layer that comprises of user and applications and lower layer comprising the operating system and communication facilities.
- This organization of distributed system is called as middleware is shown in Fig. 1.8.1.

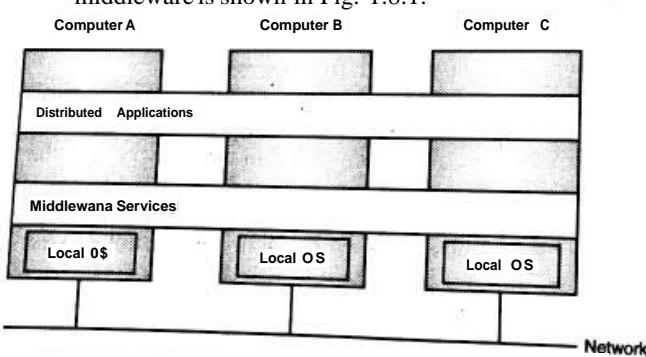


Fig. 1A1 : Distributed system organized as middleware

- In Fig. 1.8.1, applications are running on three different machines A, B and C in network.
- As distributed systems are built on top of computer networks, which are of two types LANs (Local Area Networks), and WANs (Wide Area Network).
- LAN covers one room, building or campus
- WAN covers city, country, or whole world. Ethernet is

example of LAN. Internet is a network and can be connected to WAN. Machines can communicate by exchanging packets.

Each packet contains address of source and destination. Router contains routing tables. It extracts address from incoming packet and looks up the table to find outgoing line to send the packet.

This process is repeated until the packet is delivered to the destination.

X. are very dynamic and continuously as routers and links go down and back up and as traffic conditions change.

Protocol is set of rules by which Compaq communicates with each other.

Many protocols are present such as router protocols, host-host protocols, and others. Protocols stack layers different protocols on top of one another. This protocol stack is used by all modern network. Different layer in protocol stack deal with different issues.

Most of the distributed systems use the Internet as a base. Hence, these systems use two important Internet protocols: IP (Internet protocol) and TCP (Transmission Control Protocol),

IP (Internet Protocol) is a datagram protocol. In IP, a sender sends datagram of up to 64 KB over the network and no guarantees are given for its delivery.

The datagram may be fragmented into smaller packets and travel independently, possibly along different routes.

At destination host when all packets reaches, they are

* assembled in correct order as per sequence number and delivered to the application.

IP protocol has two versions, v4 and v6. Version 4 (V4) are currently in use and v6 is up and coming. IP packet starts with a 40-byte header that contains 32 bit source and destination address with other fields.

These are called **IP addresses** and routing is carried out using these addresses.

Internet offers reliable communication in *.

Transmission Control Protocol is present on top of IP.



- TCP makes use of IP to offer connection-oriented communication on port number.
- Remote process is specified by IP address of machine and port number.
- Sender first establishes the connection and sends bytes over that connection which is guaranteed to come out in the correct order.
- The TCP gives this guarantee by using sequence numbers, checksums, and retransmissions of incorrectly received packets.
- A machine is identified by IP address and it is difficult to manage such list of huge number of IP addresses, DNS (Domain Name System) was invented as a database that maps ASCII names for hosts onto their IP addresses.
- This naming system permits the mail program on the sending machine to search the destination host's IP address in the DNS database, establish a TCP connection to the mail daemon process there, and send the message as a file.
- The user-name is sent along to identify which mailbox to put the message in.
- A tightly-coupled operating system is called as a **distributed operating system (DOS)**, and it manages multiprocessors and homogeneous multicomputer.
- The loosely-coupled **network operating system (NOS)** is used for heterogeneous multicomputer systems.
- Management of the underlying hardware is an important issue for a network operating system.
- The difference from traditional operating systems comes from the fact local services are made available to remote clients.

Syllabus Topic : Special Purpose Systems

1.9 Special Purpose Systems

Q. Explain various special purpose systems.

Apart from general purpose computer system, some classes of computer systems have limited functions and their objectives are deal with limited computation domains. These are shown in Fig. C1.3.

Types of & Special Purpose Systems

- 1. Real-Time Embedded Systems
- 2. Multimedia Systems
- 3. Handheld Systems

Fig. CL3 : Special purpose systems

→ 1.9.1 Real-Time Embedded Systems

- Embedded computers are the most common form of computers. These devices are used in every field of our life, from car engines and manufacturing robots to VCRs and microwave ovens.
- These devices perform very specific tasks. These embedded systems differ significantly.
- These can be general-purpose computers, running standard operating systems such as UNIX or Linux with special-purpose applications to implement the functionality.
- Some of the systems are hardware devices with a special-purpose embedded operating system providing just the functionality desired.
- The embedded systems are becoming more multifaceted and complex today. Also these systems will affect our life with more involvement.
- This means they will bear more and more responsibilities on their shoulders to solve real life problems to make our life easier.
- So real time operating system needs to be effective to manage more complex real time applications.
- Real time operating systems must respond quickly. These systems are used in an environment where a large number of events (generally external) must be accepted and processed in a short time.
- Real time processing necessitates quick dealing and characterized by providing instant response. For example, a measurement from a petroleum refinery indicating that temperature is getting too high and might demand for immediate attention to avoid an explosion.
- In real time operating system swapping of programs from primary to secondary is not frequent.
- Most of the time, processes remain in primary memory in order to provide quick response, therefore, memory.

Operating System (MU - Sem 4 - T)

system is less demanding management in real time compared to other system.

- The primary functions of the real are:
 - to Management of the CPU and other resources requirements of an application.
 - fulfill the inquiry.
 - o Synchronization with and responding to the system event.
 - o Efficient movement of data among processes and carrying out coordination among these processes.
- In addition to these primary functions that are now commonly included to enhance the performance. These are:
 1. To provide an efficient management of KA ..
 2. To provide an exclusive access to the computer resources.
- Few more examples of real time processing are:
 1. Airlines reservation system.
 2. Air traffic control system.
 3. Systems that provide immediate updating.
 4. **Systems that provide up to the minute information on stock prices.**
 5. Defense application systems like as RADAR.

> 19.2 Multimedia Systems

- Apart from conventional data, now days operating system should be able to handle multimedia data.
- Multimedia data comprises audio and video files as well as conventional files.
- Multimedia data should be delivered to the application in defined time constraint. E.g. Video frames must be streamed as 25 frames per second.
- Multimedia depicts a broad range of applications that are in well-liked use today.
- These comprise audio files such as MP3 DVD movies video conferencing, and short video clips of movie previews or news stories downloaded over the Internet
- Multimedia applications may also , broadcasting over the World Wide Web ^{comprise} for example, broadcasting of speeches or sports events
- Multimedia applications include combining , audio and video. For example, a movie ^{in a} ^{bosch} ^{of} separate audio and video tracks.

Multin- a 81.p. 'X" O d directed toward personal « D.p. " As M(J cellar telephones as well as devices Toneme * S.em Sho Old The design ^t TM.dia systems. also support for requirements of

> 1.9.3 Handheld Systems

Q. Write m

- **Personal Digital Assistants (PDA)** for example and pocket-PCs and cellular telephones are examples of handheld systems.
- Several of such system uses special-purpose embedded operating systems.
- These devices are having small size and have a small amount of memory, slow processors, and small display screens.
- Physical memory size of these devices is between 512 KB and 128 MB.
- Therefore it is a job of operating system and applications to manage memory efficiently.
- **A second concern to developers of handheld devices** is the speed of the processor used in the devices. Handheld devices require faster processors with compare to PC.
- Faster processors need more power. So it is obvious that, handheld devices require large battery size and will occupy more space for battery.
- Therefore most handheld devices use smaller, slower processors that consume less power.
- As a result, the operating system and applications must be designed not to tax the processor.
- The final issue deal with program designers for handheld devices is I/O. A short of physical space restricts input methods to small keyboard, handwriting recognition, or small screen-based keyboards.
- The small display screens also provide constraints on output options.
- The small display for a handheld device limits to 3 inches square with compares to 30 inches in PC.
- Therefore reading e-mail and browsing ⁱⁿ ^{8"} ^{reduced into smaller displays.}
- **A solution to this problem is to deliver small set of web page and display called as web clipping.**

Aylinbus Topic: Operating System Services**1.10 Operating System Services**

-> (May 16)

- Q. Services provided by operating system.
MU ■ May 2016. 5 Marks

Following are the six services provided by operating system for efficient execution of users application and to **make convenience to use computer system**.

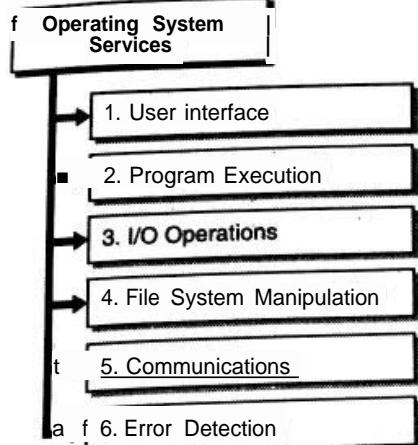


Fig. C1.4 : Operating system services

→ **1. User Interface**

- User interface is essential and all operating systems provide it. Users either interface with the operating system through command-line interface or graphical **user interface or GUI**.
- Command interpreter executes next user-specified command. A GUI offers the user a mouse-based **window and menu system as an interface**.

→ **2. Program Execution**

- The operating system provides an environment to run users programs efficiently.
- The resources needed for the programs to complete execution are provided by operating system ensuring optimum utilization of computer system.
- Memory allocation and deallocation, processor allocation and scheduling, multitasking etc functions are performed by operating system.
- The operating system has all rights of resource management. User program does not have these rights.

→ **3. I/O Operations**

- Each program requires cannot produce output without taking any input.

- This involves the use of I/O. During execution, program requires to perform I/O operations for input or output.

- Until I/O is completed, program goes in waiting state. I/O can be performed in three ways i.e. programmed I/O, Interrupts driven I/O and I/O using DMA. The underlying hardware for the I/O is hidden by operating system from users.

- These I/O operations make it convenient to execute the user program.

- Operating system provides this service to user program for efficient execution.

→ **4. File System Manipulation**

- Program takes input, processes it and produces the output. The input can be read from the files and produced output again can be written into the files.

- This service is provided by the operating system. All files are stored on secondary storage devices and manipulated in main memory,

- The user does not have to worry about secondary storage management.

- Operating system accomplishes the task of reading from files or writing into the files as per the command specified by user,

- Although the user level programs can provide these services, it is better to keep it with operating system.

- The reason behind this is that, program execution speed is fundamental to VO operations.

→ **5. Communications**

- Cooperating processes communicate with each other. If communicating processes are running on different machines then messages are exchanged among them and they get transferred through network.

- This communication can be realized by making the use of user programs by customizing it to the hardware that facilitates in transmission of the message.

- Customization of the user programs can be done by means of offering the service interface to operating system to realize communication among processes distributed system.

- In this way the communication service will be provided by operating system and user programs from will be free from taking care of communications.

→ **6. Error Detection**

- In order to prevent the entire system from malfunctioning, the operating system continually keeps watch on the system for detecting the errors.

- User programs kept free from such error detection to improve the performance.
- On the contrary, if this right would have been kept with user programs, most of the user program would have wasted time in error detection and actual work would be minimized resulting in performance degradation.
- Operating system needs to carry out complex tasks **during the process of error detection.**
- Such task comprises deal location of many resources such as processor, memory etc.
- If these tasks are again kept with user programs, then it **may disturb the normal operation of operating system.**

1.11 Objectives of Operating System

→ (June 15, Nov. 15)

Q. Explain different objectives of operating system.

MU - June 2015. Nov 2015. 4 Marks

Following are the three objectives of Operating System,

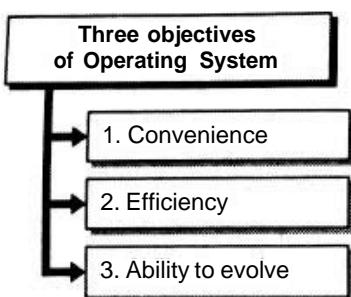


Fig. C1.5 : Objectives of operating system

→ 1. Convenience

Computer system can be conveniently used due to operating system.

→ 2. Efficiency

- Computer system comprises of many resources.
- All these resources are utilized by users application in **efficient manner due to operating system.**
- 3. Ability to evolve
- The design of the operating system permits the efficient development testing.
- It is also supports for flexibility by allowing for addition of new system functions without **interfering** with service.

1.11.1 Functions of Operating System

→ (Dec. 14, June 15, Nov. 15)

Q. Explain different functions of OS?
MU - Dec 2014. June 2015, Nov. 2015. 3 Marks

- **Operating system comprises different modules** each of having its own collection of defined inputs and outputs.
- These different modules or components of **one** system carry out specific tasks to offer the **co** functionality of the operating system.
- The most important functions of the operating system are shown in Fig. C1.6.

Functions of Operating System

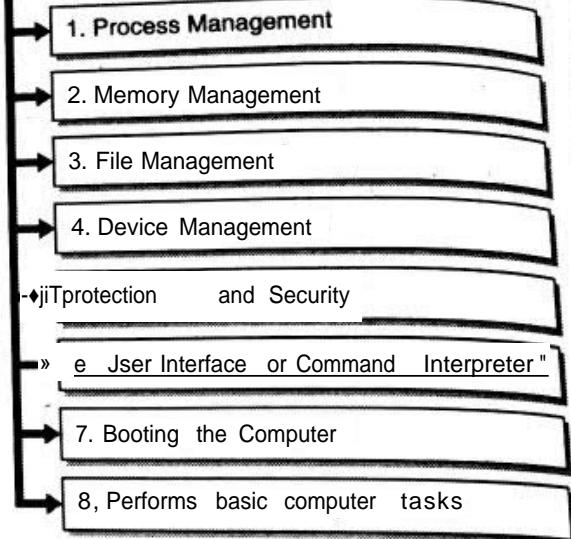


Fig. C1.6 : Functions of the operating system

→ 1. Process Management

The process management activities involves :

1. To provide control access to shared resources like file, memory, VO and CPU.
2. Control the execution of user applications.
3. Creation, execution and deletion of user system processes.
4. Resume a process execution or cancel it.
5. Scheduling of a process.
6. Synchronization, interposes communication and deadlock handling for processes.

→ 2. Memory Management

Following memory management related functions carried out by OS :



- 1 It allocates the primary memory as well as secondary memory to the user and system processes.
- 2 Reclaim the allocated memory from all the processes that have finished its execution.
- 3 Once used block becomes free, OS allocates it again to the processes.
- 4 Monitoring and Keeping track of how much memory used by the process.

→ 3. File Management

The file management activities of operating system consist of :

- 1 Files and directories are created and deleted by OS.
- 2 OS offer the service to access the files and also it allocates the storage space for files by using different methods of allocation,
- 3. It keeps back-up of files.
- 4. It offers the security for files.

→ 4. Device Management

The device management tasks include :

- 1. Device drivers are opened, closed and written by OS.
- 2. Keep an eye on device driver. Communicate, control and monitor the device driver.

→ 5. Protection and Security

- The resources of the system are protected by the operating system.
- In order to offer the needed protection, the operating system such as read, write, encryption, and back-up.

→ 6. User Interface or Command Interpreter

User interacts with computer system through operating system. Hence OS act as an interface between the user and the computer hardware. This user interface offers

through set of commands or a graphical user interface (GUI). Through this interface user can interact with the applications and the machine hardware.

4 7 Booting the Computer

- The process of starting or restarting the computer is known as booting.
- If computer is switched off completely, then it is called cold booting. Booting is the process of using the operating system to restart the computer.

→ 8. Performs basic computer tasks

- The management of various peripheral devices such as the mouse, keyboard and printers is carried out by operating system.
- Today most of the operating systems are plug and play. These operating systems automatically recognize and configure the devices with no user interference.

Syllabus Topic : Operating-System Interface

1.12 User Operating-System Interface

| Q. Write short note on user-OS interface.

>User Interface

- User interface is essential and all operating systems provide it.
- Users either interface with the operating system through command-line interface OR graphical user interface or GUI
- Command interpreter executes next user-specified command.
- A GUI offers the user a mouse-based window and menu system as an interface.

Graphical User Interface or Command Interpreter

- User interacts with computer system through operating system. Hence OS act as an interface between the user and the computer hardware.
- This user interface offered through set of commands or a Graphical User Interface (GUI).
- Through this interface user makes interaction with the applications and the machine hardware.
- GUI stands for graphical user interface. It permits user to use icons or other visual indicators to interact with computer, rather than using only text via the command line.
- Due to GUI, it is not necessary to remember the commands. It provides the ease of use interaction with computer. user need to know any programming language for interacting with computer.

Linux, Apple, Mac OS and some of the systems that offer GUL

Syllabus Topic : System Calls**1.13 System Calls**

→ (June 15, Nov. 15, May 16)

- Q. What are system calls?
MU - June 2015, Nov. 2015. 4 Marks
- O. Write short note on system calls.
MU - May 2016. 5 Marks

- The interface between OS and user programs is defined by the set of system calls that the operating system offers. **System call is the call for the operating system to perform some task on behalf of the user's program.**
- Therefore system calls make up the interface between processes and the operating system.
- The system calls are functions used in the kernel itself. **From programmer's point of view, the system call is a normal C function call.**
- Due to system call, the code is executed in the kernel so that there must be a mechanism to change the process mode from user mode to kernel mode.
- In the UNIX operating system, user applications do not have direct access to the computer hardware.
- Applications first request to the kernel to get hardware access and access to computer resources.
- During execution when application invokes a system call, it is interrupted and the system switches to kernel space.
- The kernel then saves the process execution context of **interrupted process and determines purpose of the call.**
- The kernel safely makes sure that the request is valid and that the process invoking the system calls has **sufficient privilege.**
- **Some of the system calls can only be invoked by a user with super user privilege.**
- If the whole thing is fine, the request in kernel mode and can access the drivers in charge of controlling the hardware
- The data of the calling process can be modified by kernel, as it has access to memory space. But it cannot execute any code in user application, for clear security reasons.
- When the kernel finishes the request, it stores the process execution when the system call was invoked, that is saved

After this, control returns to the calling program which persists executing.

USER SPACE : Application and Libraries**KERNEL : Process, memory, device, and file management etc.**

~ hardware"

TM dav .. se

7

Fig. 1.13.1: The kernel

- As shown in Fig. 1-13-1 the kernel is a central module of most computer operating systems. Its main responsibilities are to be shown in Fig. C1.7.

Main responsibilities of Kernel

- 1. Act as a standard interface to the system hardware
- 2. Manage computer resources
- 3. Put into effect isolation between processes
- 4. Implement multitasking

Fig. C1-7: Responsibilities of kernel

→ **1. Act as a standard interface to the system hardware**

For example, while reading a file, application does not have to be aware of the hard-drive model or physical geometry as kernel provides abstraction layer to the hardware.

→ **2. Manage computer resources**

As several users and programs share machine and devices, access to those resources must be synchronized.

The kernel implements and ensures a fair access to resources such as the processor, memory and the disk.

→ **3. Put into effect isolation between processes**

The kernel guarantees the execution environment of another process. Any process cannot access the memory that has been allocated to other process.

→ **4. Implement multitasking**

Process gets the false:

- Actually, several processes use constant "f", the kernel keeps switching the processor behind the scene

Ways for the process to switch from the user mode to the kernel mode

There are two ways for process to switch from the user mode to the kernel mode. These are:

- A user process can explicitly request to enter in kernel mode by issuing a system call.
- During the execution of user process kernel can take over to carry out some system housekeeping task.
- The kernel mode is both a software and hardware state. Modern processors offer a advanced execution mode, called as Supervisor Mode in which only kernel runs.
- The privileged operations are such as modifying special registers, disabling interrupts, accessing memory management hardware or computer peripherals.
- If it is not in supervisor mode, the processor will reject these operations.
- System calls take place in different ways, depending on the computer in use. Apart from the identity of the desired system call, more information is needed.
- The precise type and amount of information differ according to the particular operating system and call. Consider the example of getting the input.
- We may require specifying the source, which can be file or device.
- We also need to specify the address and length of the memory buffer into which the input should be read.
- The device or file and length may be implicit in the call.

Ways of parameter passing to operating system

Parameters can be passed to the operating system by following three different ways:

- Pass the parameters in registers.
- If there are more parameters than registers, then the parameters are generally stored in a block, or table, in memory, and the address of the block is passed as a parameter in a register. This is the approach taken by Linux and Solaris.
- Program can place the parameters onto the stack which then popped off the stack by the operating system.

Some operating systems favor the block or stack method, because those methods do not limit the number or length of parameters being passed.

Syllabus Topic : Types of System Calls

1.13.1 Types of System Calls

→ (June 15, Nov. 15)

- Q. Explain any five system calls.

MU - June 2015. Nov. 2015. 6 Marks

System calls are categorized in five groups. These are shown in Fig. C1.8.

Types of System Calls

1. Process control
2. Device manipulation
3. Communications
4. File manipulation
5. Information maintenance

Fig. C1.8 : Types of system calls

Sr. No.	Group	Examples
1.	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
2.	Device manipulation	request device, release device, read, write, reposition, get device attributes, set device attributes, logically attach or detach devices
3.	Communications	create, delete communication connection, send, receive messages, transfer status information, attach or detach remote devices
4.	File manipulation	create file, delete file, open, close, read, write, reposition, get file attributes, set file attributes



Sr. No.	Group	Examples
5.	Information maintenance	get time or date, set time or date, get system data, set system data, get process, file, or device attributes. set process, file, or device attributes

1.13.2 Some Examples of System Calls

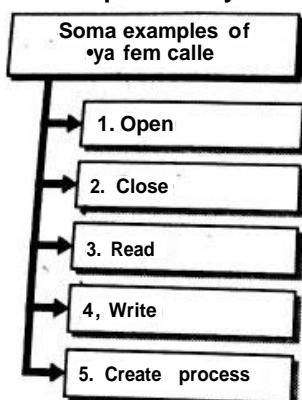


Fig. Cl.» : Examples of system calls

→ 1. Open

- Open system call request to the operating system for using a file.
- The file's path is used as argument to specify the file.
- The 'flags' and 'mode' arguments to this call specifies how to use the file.
- On successful approval by operating system, it returns a "file descriptor" which is a positive integer.
- **Returning of value -1 indicates denial of access to file** and "errno" needs to be checked to get reason of denial.

* X Close

- **Close()** informs to the operating system that file with said file descriptor is used.
- After this, same file descriptor can be reused by operating system.

→ 3. Read

- **Read()** specifies to the OS the number of bytes to read from the file opened in file descriptor "2".
- It also specifies the location to which the bytes are to be written.
- The return value is number of bytes actually read.

→ 4. Write

- **Write()** is similar to **read()** system call, only it writes the bytes instead of reading them.
- It returns the number of bytes actually written, which is almost invariably "size".
- 5. Create process
- **During the execution process can create new process** using create process system call.
- The process which creates new process is parent process, and the new processes are children of that process.
- Newly created processes may in turn create processes, creating a tree of processes.

Syllabus Topic : System Programs

1.14 System Programs

- Q. Explain various system programs that are associated with operating system.

- System programs offer a suitable environment for development of application programs and its execution.
- Out of these system programs, a number of are on user interfaces to system calls and rest of the system programs are much more complex.
- Operating systems are supplied with these programs and are classified as below :

0 **File management** : System programs in this category usually are responsible for manipulation of files and directories. These programs generally create, delete, copy, rename, print, dump and list the files.

0 **Status information** : The status information

the system is for example: date, time, total free memory or disk space in hand, number of used Present in the system, etc.

Some system programs simply inquire for information. Other gives the information regarding complete performance, logging, and debugging.

System programs have more complex characteristics, these programs carries formatting and printing of output to the terminal.

- files or other output devices or display it in a window of the graphical user interface.
- o **File modification** : System program such as text editor is used to create and modify data or information in files stored on disk or other storage devices. Some particular commands are used to look for contents of files or carry out conversions of the text.
 - o **programming-language support** : The language translators such as compilers, assemblers, interpreters come with operating system to translate the user programs to object code.
 - o **Program loading and execution** : After translation of the user program in object code it needs to be loaded into memory for carrying out the execution. The different types of loaders such as absolute loaders, relocatable loaders, linkage editors, and overlay loaders are offered by the system in order to load the program in memory.
- Systems required for debugging the higher-level languages or object code are also offered.
- o **Communications** : Some of the system programs offer the way for creating virtual connections among processes, users, and computer systems. They permit users to exchange messages, to browse web pages, to send electronic-mail messages, remote login and transferring the files from one machine to another.
 - Just like the system programs, many applications too are the part of operating system.
 - These programs are called as system utilities or **application programs**.
 - It comprises the applications such as web browsers, word processors and text formatters, spreadsheets, database systems, compilers, plotting and statistical-analysis packages, and games.

1.14.1 Comparison between System Program and Application Program

Q. Compare system and application program.

Sr. No.	Parameters	System Program/System Software	Application Program/Software
1.	Definition	System programs are designed to manage and control computer hardware to help the users program for efficient execution and to offer the services to them .	These software helps the user to carry out a specific tasks as per users requirement.
2.	Purpose	It is general-purpose.	Its purpose is <u>specific</u> .
3.	Environment	System software is itsen responsible to create environment to execute itself and other application programs.	Environment required for application program is created by system software.
4.	Responsibility	It is responsible to control and manage entire computer system.	It controls particular task for which it was intentionally designed.
5.	Execution	It executes all the time when system is on to offer the services to user programs.	Its execution is as per the need of user.
6.	Examples	Operating system, compilers, loaders, text editors. assemblers are some of the examples of system programs	Word processors or any software designed to perform specific task.

Syllabus Topic : Operating System Design and Implementation

1.15 Operating System Design and Implementation

There are many problems in the design and implementation of operating systems. There are certain solutions to these problems.

1.15.1 Design Goal*

Q. Explain design goals of operating system

is to define

- The initial problem that needs to be taken is to define goals and specifications.
- Which hardware is considered and what type of system is decided to design will affect the design of the system.
- Types of system include the single or multiuser, real time, distributed, general purpose, interactive or batch system. This is the highest level of design.
- It is very difficult to identify the requirements. These requirements are categorized in two fundamental groups: *user goals* and *system goals*.
- **Users always wishes the system which is convenient to use, simple to learn and to make use of. It should be safe and protected and having high performance.**
- Due to lack of general agreement on how to accomplish these goals, the above properties are not useful in the design of system.
- Person responsible to design, create, maintain, and operate the system also desires the properties like: Simple to design, implement, and maintain, flexibility, reliability, error free, and efficient. Such requirements are unclear and may be inferred in different ways.
- Single solution to the problem of defining the requirements for an operating system does not exist.
- The different systems in existence show that different requirements can lead to a large diversity of solutions for different environments.

1.15.2 Separating Policies from Mechanisms

- Mechanisms decide the way of doing something whereas policies determine what will be done.
- It is necessary to keep policies separate from mechanisms to achieve the flexibility

- Changes in policies occur over the period of time. At extreme level, each change in policy would need change in the underlying mechanism.
- A mechanism which does not affect by changes in policy would be more preferred.
- If policy is changed then certain parameters of mechanism are again

Another example of the mechanism for priority to certain types of jobs

If "Xel'y" mechanism can be used to support that I/O-bound jobs should priority over CPU-bound jobs or to support the reverse.

The essential set of primitives having microkernel design.

Consequently, it allows to add more mechanisms and policies by means of user kernel modules or by means of user programs themselves.

The allocation of resources should be on the basis of policy decision.

1.15.3 Implementation

Q. Explain implementation of operating system design.

- Traditionally assembly language was used to implement the operating system.
- Now high level languages are used to implement the operating system.
- Due to this code is easy to read, understand and development and implementation process became faster.
- After recompilation operating system can be ported to other machine.
- On the other side, if high level languages are used to implement the operating system the speed will be reduced and storage required would increase.
- These issues are negligible as assembly W routines can be developed for large programs and compilers can optimize the code.
- The operating system routines responsible for bottleneck can be replaced by assembly W routines after system functions correctly.
- In this way performance can be improved. Use of data structures and algorithms also improves performance.
- **Operating system should also support for monitoring system performance.**

- Bottlenecks can be identified by monitoring system performance.
- A static system should include the code to calculate Xdisplay measures of system behaviour.
- In many systems, the operating system carry out this job by producing trace listings of system behaviour.
- All necessary events are logged with their time and important parameters are written to a file.

Syllabus Topic: Virtual Machines

1.16 Virtual Machines

q. What is virtual machine?

- By means of CPU scheduling and virtual machine (VM) techniques, host OS can produce illusion that each process has its own processor and memory.
- This memory is virtual memory. The VM offers an interface that is similar to the underlying bare hardware.
- The processes running are the guest processes and each process gets virtual copy of the underlying computer. This guest process is an operating system.
- In this way single machine runs the multiple operating systems in parallel, each in its own VM.
- From availability and security point of view, using separate computer to put each service is more preferable by organizations. If one server fails, other will not be affected.
- It is also beneficial in case if organizations need to use different types of operating system.
- However, keeping separate machines for each service is costly. Virtual machine technology can solve this problem. Although this technology seems to be modern but idea behind it is old.
- The VMM (Virtual Machine Monitor) creates the illusion of multiple (virtual) machines on the same physical hardware.
- A VMM is also called as a hypervisor. Using virtual machine, it is possible to run legacy applications on operating systems no longer supported or which do not work on current hardware.
- Virtual machines permit to run at the same time applications that use different operating systems.
- Several operating systems can run on single machine without installing them in separate partitions.
- Virtualization technology plays major role in cloud computing.

1.16.1 History

- Internet standards virtual machines were also present earlier. Around 1960, IBM developed two hypervisors SIMMON and CP-40 which was research project.
- CP-67 was reimplemented version of CP-40. CP-67 formed the CP/CMS which was virtual machine OS for the IBM System/360 Model 67. In 1972, CP-67's reimplemented version was launched as VM/370 for the System/370 series.
- In 1990 IBM replaced System/370 by System/390. In all this journey, although machines were renamed but underlying architecture remained unchanged to maintain backward compatibility.
- There was better improvement in hardware technology and newer machines were bigger and faster as well.
- Over the period of time, although hardware technology improved, but virtualization was there and supported by all the machines.
- In 2000, IBM released the z-series which was backward compatible with the System/360. Z-series had supported 64-bit virtual address spaces.
- All of these systems supported virtualization decades before it became popular on the x86.
- The early releases of OS/360 were strictly batch systems. But, due to requirement, many 360 users, decided to have timesharing, so various groups, both inside and outside IBM decided to write timesharing systems for it.
- The TSS/360, the first time sharing system, arrived late and it was so big and slow that few sites converted to it. It was finally neglected. And its development cost was very high around some \$50 million (Graham, 1970).
- But a group at IBM's Scientific Center in Cambridge, Massachusetts, produced a radically different system that IBM eventually accepted as a product, and which is now widely used on its remaining mainframes.
- This system, originally called CP/CMS and later renamed VM/370 (Seawright and MacKinnon, 1979), was based on an astute observation :
 - o A timesharing system provides
 - o Multiprogramming and
 - o An extended machine with a more convenient interface than the bare hardware.
- The essence of VM/370 is to completely separate these two functions.

- The heart of the system, known as the virtu-1 machine monitor, runs on the hardware multiprogramming, providing not one, but several virtual machines to run up, ds

Fig. 1.16.1,

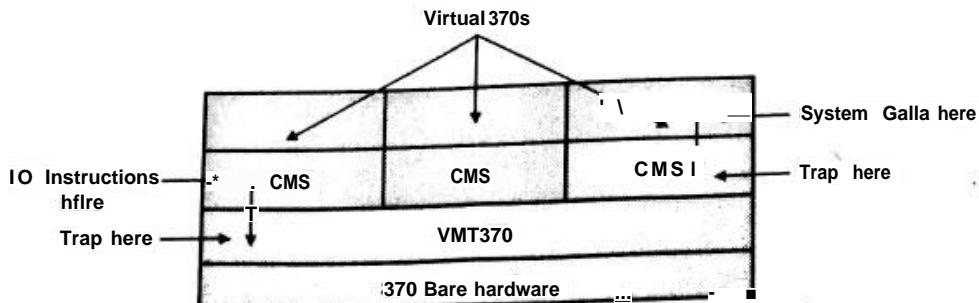


Fig. 1.16.1

- However, unlike all other operating systems, these virtual machines are not extended machines, with files and other nice features.
- Instead, they are exact copies of the bare hardware, including kernel/user mode, I/O, interrupts, and everything else the real machine has.
- To maintain the backward compatibility, the defects in Intel 386 were automatically carried forward into new CPUs for 20 years.
- Hence virtualization has been a problem on x86 architecture.
- The same instructions behave differently when executed in kernel mode with compare to their execution in user mode, for example, instructions carrying out I/O and changing MMU settings.
- All these instructions are called as sensitive instructions. Some other instructions cause trap if run in user mode called as privileged instructions.
- If sensitive instructions are subset of privileged instructions then machine is virtualizable.
- In 2005, Intel and AMD added virtualization in their CPUs. On the Intel CPUs it is called **VT** (Virtualization Technology); on the AMD CPUs it is called **SVM** (Secure Virtual Machine).
- Containers are created in order to run virtual machines in it. When a guest operating system is initiated in a container, it carries on running there until it causes exception and traps to the hypervisor, for example by executing an I/O instruction.**
- The set of operations that trap is controlled by a hardware bitmap set by the hypervisor.

1.16.2 Benefits

- Q. What are the benefits of virtual machine?

Following are the benefits of creating virtual machines.

- Several virtual machines are protected from each other and also host system is protected from viruses.
- If guest operating system contains virus then it damage this operating system but not other operating system. Also host operating system is not affected by this virus.
- The direct sharing of resources is not present. This sharing is possible.
- The network of VMs is possible where sharing can be done over virtual communication network. The research and development in operating system is possible due to virtual machine concept.
- System remains unavailable till changes made are tested in operating system.
- Time required for this is system development time: in case of virtual machine, system development is done on virtual machine instead of on a physical machine.
- The virtualized workstation permits for quick prototyping and testing of programs in changeable environment.
- Quality-assurance engineers can test applications in several environments with varying power and maintaining a computer for environment.
- VMware by combining two separate physical machines is possible by creating

- In this manner, two lightly loaded system can be converted in one heavily loaded system.
- For extensive acceptance, the design of virtual machines must be standardized with the intention that any virtual machine will run on any virtualization Platform.
- **The "Open Virtual Machine Format" can be realized to be successful in combining virtual-machine formats.**

1.16.3 Simulation

- Because of virtualization, guest OS and applications believes to be executing on native hardware.
- Simulation is other system emulation methodology like virtualization in which host system has one architecture and compilation of guest system was carried out on other architecture.
- This is just like running instructions on new computer system that were compiled on the old computer system.
- The programs could be executed in an emulalor that converts each of the old systems instructions into the native instruction set of the new computer system.

1.16.4 Para-Virtualization

- Para-virtualization offers to the guest similar but not identical preferred system of the guest.
- It is necessary to modify the guest to run on Para-virtualized hardware.
- Due to this resources can be utilized more efficiently. A virtualization layer is required. Solans 0 operating system contains zones or containers produce virtual layer between OS and applications.
- In this case, hardware virtualization is not done but and associated devices are virtualized so that processes on the within containers believe that only they are system.
- Containers can be more than one each having its own applications, network stacks, network address ports user accounts, and so on.
- CPU resources are partitioned among the containers and system wide processes.

1.16.5 Implementation

- **The implementation of virtual machine (VM) concept is difficult.**

- The creation of exact copy of underlying machine needs most of the work,
- This underlying machine has user mode and kernel mode, Since VM software is OS, it can run in kernel mode.
- The VM itself can run in user mode. It is necessary to have virtual user mode and virtual kernel mode.
- Both of these modes must run in physical user mode. Just like user mode to kernel mode transition on real machine, there should be transition from virtual user mode to virtual kernel mode on VM.

This transfer can be carried out as follows :

- System call by program running on VM in virtual user mode causes transfer to the VM monitor in the real machine.
- VM monitor then alter the register content and program counter for VM to simulate the effect of system call.
- The VM is the restarted which is now in virtual kernel mode. Virtual I/O may take less or more time than real I/O. As CPU is being multiprogrammed among many VM*s, they can be further slows down. For virtualization, hardware support is also important. All general purpose CPU supports for virtualization.

1.16.6 Examples

Following two, VMware workstation and java virtual machine are the examples of VMs.

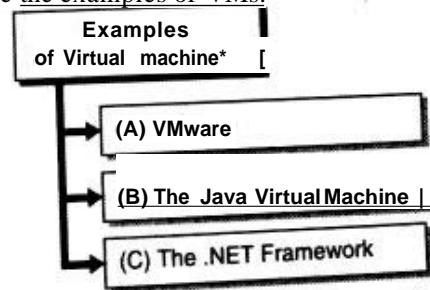


Fig. CLIO : Examples of virtual machine

-> 1.16.6(A) VMware

Explain VMware in detail.

VMware workstation is a established marketable conceptualized Intel X86 and SXhXJ. VMware Workstation runs as an application on a host machine.

OS Examples of host OS are Widows or Linux.

- This permits this host system to run a number of different guest operating systems as **independent virtual machines**.

Appikalton

- In Fig. 1.162. Linux is running as **host OS**. Three operating systems, free **BSD**, **Windows NT** and **Windows XP** are running as **guest operating systems**.

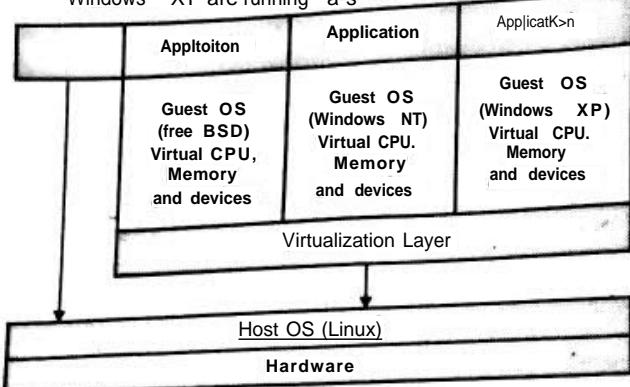


Fig. 1.16.2: VMware Architecture

- The virtualization layer offer abstraction of physical hardware into isolated VMs which are running as guest operating systems.
- As shown in Fig. 1.16.2, each VM has its own virtual memory, CPLF, devices and network interfaces etc.
- The guest file is the copy of host OS file in file system.
- The guest instance of file is separate copied file due to which protection of guest instant against disaster is achieved.

→ 1.16.6(B) The Java Virtual Machine

Q. Write short note on JVM.

- Java is object oriented programming language in which program contains one or more classes.
- The compiler produces architecture independent bytecode for each class file.
- This bytecode can run on any implementation of the java virtual machine.
- The JVM has class loader and java interpreter. The java interpreter runs the architecture independent bytecodes
- **The compiled .class files are loaded by class loader from both the Java program and the Java API which are then executed by java interpreter.**
- Once the loading of the class is done it **checks whether it is underflow or overflow the stack.**

It also makes sure that the bytecode does not carry on pointer arithmetic, which could lead to illegal memory access.

- If this verification is successful then it is run by the Java interpreter. The JVM also proves collection automatically.
- The implementation of JVM can be carried out on the top of host OS or as a part of browser.
- The interpreter interprets bytecode one at a time. faster interpretation is achieved by just in compiler.
- **The JIT first converts bytecode in machine instructions and subsequent invocation of the methods are carried out using these cached machine instructions.**
- It avoids bytecode interpretation in subsequent invocation.
- **The hardware implementation of JVM can be carried out in chip which is faster than software implementation.**

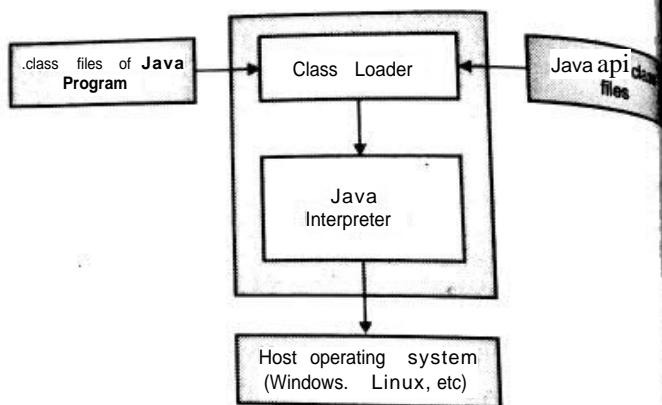


Fig. 1.163 : The java virtual machine

→ 1.16.6(C) The .NET Framework

- The .NET Framework includes set of class libraries and executing environment which creates platform for developing the software.
- A program written for .NET framework does not require specific hardware or operating system.
- The execution of the program can be carried successfully by any architecture which implements .NET.
- **The developer required for this execution is abstracted by .NET framework and VM is provided as interface between execution environment and underlying architecture.**

- The .NET framework contains common language runtime (CLR) which is implementation of .NET VM.
- The programs are written in C# and VB.NET and compiled into Microsoft Intermediate Language (MS-IL) which is intermediate architecture independent language.
- The compiled files (assemblies) having extensions as .EXE or DLL contain MS-IL instructions and metadata.
- The CLR loads assemblies into the Application Domain upon program execution.
- As executing program requires instructions, the CLR translates the MS-IL instructions inside the assemblies into native code that is precise to the underlying architecture using JIT compilation.
- The instructions then will carry on running as native code for the CPU.
- The architecture of the CLR for the .NET framework is shown in Fig. 1.16.4.

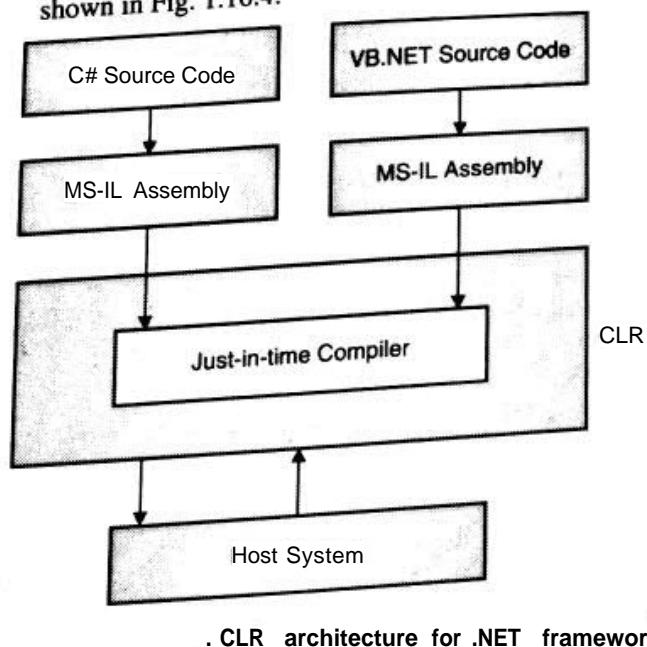


Fig. 1.16.4

TOPIC: Operation System Debugging

1.17 operating System Debugging

Q. Explain OS debugging in detail.

Debugging includes finding and fixing the errors in system. Debugging activity is carried out for both hardware and software. By fixing the bugs in system, performance can be improved,

1.17.1 Failure Analysis

- Operating system writes error information in log file or takes core dump in case of failure of process.
- The core image is maintained in file for afterward analysis.
- Debugger tool probes the executing programs or core dump. It is challenge to debug the user level process.
- As kernel is more complex and having large size, its debugging is complex task as user level debugging tools are not available. When kernel crashes, its memory state is saved to crash dump.
- The tools for OS debugging are different than for process debugging.
- If failure occurs in file system code, then it is possible for kernel to write log in file in file system before rebooting. In this case, memory state is saved on the disk area where file system is not present.
- After rebooting, this information is gathered and saved in file of file system for later analysis.

1.17.2 Performance Tuning

- It is necessary to monitor the system performance to discover bottlenecks.
- It is also important to add the code for computing and displaying the measures of system behavior.
- Operating system produces the listing of system behavior. A log is maintained for analysis with their runtime and important parameters.
- This log is then written to a file. Afterward, processing is carried out by an analysis program to know system performance and to recognize bottlenecks and inefficiencies. Recommended improved system of a
- The simulation can be carried out with these same traces as input errors in Traces also can facilitate people to locate operating-system behavior.
- In other method of performance tuning, a system can have interactive tools which permits the users and administrators to ask the different components state of the system to find bottlenecks.

1.17.3 DTrace

- DTrace is the facility that permits to add probes dynamically in running user processes and kernel probes.

Operating System (MU* Sem 4 - IT)

- By giving queryTM language. surpd«"« am system sure, and proves. ~~the impossible to debug interaction between user level~~
 - ~~about hc kernel, irjc5 can be determined,~~
 - ~~the codes and~~
 - ~~the toolset should be to debug all the area of the system and show a major impact as per use.~~
 - ~~ZZdl -se requirement, and one safe, dynamic, low impact debugging environment.~~

'Syllabus Topic : OporatIng System GeneraUon

1.18 Operating System Generation

- o The available devices, their device number > 1, device type and model of device.
 - o mTeXpt number and characteristics.
 - o The preferred OS options, parameters
 - o used The number of buffers and size to by values, the required CPU scheduling a gorithm the maximum number of processes supported and so on.
 - A system administrator can use this information to modify os source code.
 - The compilation of modified OS is then carried out to produce output-object version as per system description.
 - Another approach, tables are created and modified selected from precompiled library.
 - After linking the modules, generated OS is formed
 - next approach, complete table driven system created.
 - Code selection is at execution time rather than compilation or link time*

Syllabus Topic : System Boot

1.19 System Boot

- When system starts up the RAM initially is in unknown state. ROM is suitable because it needs no initialization and cannot be infected by a computer virus.
- The bootstrap program runs diagnostics to decide the state of the machine. If the diagnostics pass, the program can carry on with the booting steps.
- It also initializes all parts of the system, from CPU **registers to device controllers and the contents of main memory**.
- It then starts the operating system. Several systems like cellular phones, PDAs, and game consoles store the complete operating system in ROM.
- If the operating systems are small in size, storing it in ROM is appropriate for, simple supporting hardware, and rugged operation.
- A difficulty with this approach is that changing the bootstrap code needs changing the ROM hardware chips.
- This problem is solved by some system using **Erasable Programmable Read-Only Memory (EPROM)**.
- EPROM is read only apart from when clearly given a command to become writable.
- The characteristics of ROM go down somewhere between those of hardware and those of software. Therefore all forms of ROM are also known as **firmware**.
- For large size operating systems like Windows, Mac OS X and UNIX or for systems that change often, the bootstrap loader is stored in ROM, and the operating system is on disk.
- In this case, the bootstrap runs diagnostics and has a small piece of code that can read a single block at a fixed location (say block zero) from disk into memory and execute the code from that boot block.

1.20 ExamPack (University and Review Questions)—

↳ Syllabus Topic : Introduction

Q. What is operating system?

(Refer section 1.1) {2 Marks}

(Dec. 2014, June 2015, Nov. 2015)

or Syllabus Topic : Operating System Structure

Q. Explain monolithic system. (Refer section 1.2.1)

- Q. Explain layered system in detail.
(Refer section 1.2.2)
- Q. Explain client-server model in detail.
(Refer section 1.2.3)
- Q. Differentiate between monolithic and microkernel.
(Refer section 1.2.4) (5 Marks)
- (June 2015, Nov 2015)
- Q. What is Kernel ? Describe briefly the approaches of designing Kernel.
(Refer section 1.2.4) (5 Marks)
- (Dec. 2016)
- Q. Differentiate between monolithic and microkernel.
(Refer section 1.2.4) (5 Marks)
- (June 2015, Nov. 2015)

☞ Syllabus Topic : Operating System Operations

- Q. Explain operating system operations.
(Refer section 1.3)

tr Syllabus Topic : Process Management

- Q. Explain process management function of operating system. (Refer section 1.4)

☞ Syllabus Topic : Memory Management

- Q. Explain memory management function of operating system. (Refer section 1.5)

& Syllabus Topic : Storage Management

- Q. Explain file management in OS.
(Refer section 1.6.1)

- Q. Write note on mass storage management.
(Refer section 1.6.2)

- Q. Explain how caching is implemented ?
(Refer section 1.6.3)

ur Syllabus Topic : Protection and Security

- Q. Explain protection and security mechanism of OS.
(Refer section 1.7)

cF Syllabus Topic : Distributed System

- Q. What are the characteristics of distributed system?
Explain. (Refer section 1.8)

☞ Syllabus Topic : Special Purpose Systems

- Q. Explain various special purpose systems.
(Refer section 1.9)

ffrrwetinoSystemfW:

- Q. write note on handheld syrrfe-™.
(Refer section 1.9.3)
- Syllabus Topic : operating System Services**
- Q. Explain services provided by **(May 2016)**
(Refer section 1.10) (5 Marks)
- Q. Explain different objectives of operating system.
(Refer section 1.11) (4 Marks)
(June 2015, Nov 2015)
- Q. Explain different functions of OS?
(Refer section 1.11.1) (3 Marks)
(Dec 2014, June 2015, Nov. 2015)

☛ Syllabus Topic : Operating-System Interface

- Q. Write short note on user-OS interface.
(Refer section 1.12)

☛ Syllabus Topic : System Calls

- Q. What are system calls ?
(Refer section 1.13) (4 Marks)
(June 2015, Nov. 2015)
- Q. Write short note on system calls.
(Refer section 1.13) (5 Marks) ("ay 2016)

☛ Syllabus Topic : Types of System Calls

- Q. Explain any five system calls.
(Refer section 1.13.1) (6 Marks)
(June 2015, Nov. 2015)

☛ Syllabus Topic : System Programs

- Q. Explain various system programs that are associated with operating system.
(Refer section 1.14)

Comp are system and application program.

(Refer section 1.14.1)

☛ Syllabus Topic : Operating System Design Implementation

- a. Explain design goals Of Operating system.
(Refer section 1.15.1)

- Q. Explain implementation of operating system design
(Refer section 1.15.3)

☛ Syllabus Topic : Virtual Machines

- Q. What is virtual machine? (Refer section 1.16)

- Q. What are the benefits of virtual machine?
(Refer section 1.16.2)

- Q. Explain VMware in detail. (Refer section 1.16.6(A))

- Q. Write short note on JVM. (Refer section 1.16.8J)

☛ Syllabus Topic : Operating System Debugging

- Q. Explain OS debugging in detail. (Refer section 1.17)

☛ Syllabus Topic : Operating System General

- Q. Explain operating system generation.
(Refer section 1.18)

☛ Syllabus Topic : System Boot

- Q. Explain booting process of the system in detail.
(Refer section 1.19)

CHAPTER

2

Module II

Process Management

Syllabus Topics

Process concept : Process Scheduling, Operation on process and Interprocess communication; Multithreading, Process : Multithreading models and thread libraries, threading issues; Process Scheduling; Basic concepts, Scheduling algorithms and Criteria, Thread Scheduling and Multiple Processor Scheduling,

2.1 Introduction

- Process manager implements the process management functions. In multiprogramming, single CPU is shared among many processes. If many processes remain busy in completing I/O, CPU is allocated to only one process at a given point of time.
- Here some policy is required to allocate CPU to process, called as CPU scheduling. If multiple users are working on the system, operating system switches the CPU from one user process to other.
- User gets the illusion that only he or she is using the system. Process synchronization mechanism is required to ensure that only one process should use critical section.
- Process communication, deadlock handling, suspension and resumption of processes and creation and deletion of the processes etc are some of the activities performed in process management.

Syllabus Topic : Process Concept

2.2 Process Concept

→ (Dec. 14)

q. What do you mean by process?

MU - Dec. 2014. 3 Marks

Q. Explain the concept of process. *

A program or application under execution is called as process. A process includes the execution context. A

program resides on the disk. On disk it does not require any resources.

- A program gets executed in main memory. So it should be transferred from disk to main memory.
- To complete execution, program needs many resources and competes for it.
- Now it becomes process. From the computation context point of view, a process is defined by CPU state, memory contents and execution environment.
- A CPU state is defined by the content of the various registers such as Instruction Register (IR), Program Counter (PC), Stack Pointer (SP) and general purpose registers.
- A small amount of data is stored in CPU registers. Memory contains program code and its predefined data structures.
- Heap is reserved memory area for dynamically allocation of memory to the program at run time.
- In stack, program local variables are allocated and return values of function call are stored. Some register values are also saved in stack.
- Execution environment includes open files, communication channels to other processes etc. following are the components of the process :
 - o The object code that is to be executed.
 - o Resources required by the program to complete the execution.
 - o The data on which program will operate.
 - o Program execution state.

2.3 Context Switch

- Q. What is context switch?
- When CPU switches from one process to another, a context switch occurs. The CPU (Central Processing Unit) from one process or thread to another.
 - When context switch occurs, the information of currently executing process is saved to the same later use. When CPU needs to resume the process, the restored information is used to resume the execution.
 - The context of a process is represented by Control Block (PCB) of a process.
 - The information needs to be restored to execute the next instruction to be executed (program counter), CPU register contents, pointers to memory allocated to the process, scheduling information, changed state, I/O state, etc.
 - While context switch, system does not perform any useful work. So context switch is pure overhead on the system.
 - The speed of context switching depends on the number of registers that must be copied, memory speed, etc. Context switch speed varies from system to system.

Syllabus Topic : Operation on Processes

2.4 Operations on Processes

- Q. What are the operations performed on process?

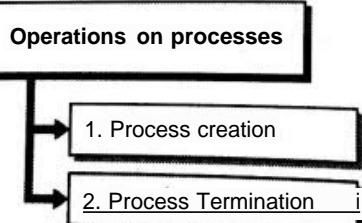


Fig. C2.1 : Operations on processes

2.4.1 Process Creation

- Q. Explain process creation.

Processes are created because of the following principal events. These are shown in Fig. C2.2.

Four principal events of process creation

1. system initialization
2. process runs
3. A user request to create a new process
4. Starting of a batch job

Fig. C2.2 = Events of process creation

System tasks

- After booting OS creates many processes.

- If running process calls

Currently executing process creates recesses by issuing create process system call. It needs these processes to assist it in execution.

Creating new processes is mainly helpful work to be carried out can easily be formulated of several related, but otherwise interacting processes.

3. A user request to create a new process

- In interactive systems, users can start a process by typing a command or (double) clicking an icon.

4. Starting of a batch job

- In this case users can submit batch jobs to the remote system. When the OS makes a decision to execute another job, it creates a new process to execute the next job from the queue in it.

During the execution process can create new processes using *create process* system call.

- The process which creates new processes is called **parent** process, and the new processes are called **children** of that process.

- Newly created processes may in turn create new processes, creating a **tree** of processes.

- In UNIX or the Windows family of operating systems, processes are identified by unique Process ID (or pid), which is typically an integer. The command is used in UNIX to obtain the process ID.

- By using command `ps` - we can obtain information for all processes currently active in the system.
- Process requires certain resources like CPU time, memory, files, I/O devices to complete its task.
- The subprocess also needs these resources. Subprocess can get its needed resources directly from the operating system, or it may be forced to a subset of the resources of the parent process.
- The parent can share the resources among its children or it can divide the resources to allocate to its children.
- If child process is restricted to a subset of the parent's resources overloading the system by creating too many subprocesses can be avoided.
- The initialization data may be passed along by the parent process to the child process after creation of the process.

After creation of a new process, two possibilities exist related to execution :

- The parent and its children both execute in parallel.
- Until few or all the children terminate, parent will wait and will not be terminated.

With respect to the address space of the new process again two possibilities exist :

- Program and data of child and parent process is same. It means child is duplicate of its parent.
- Program and data of child and parent process is different. It means child process has a new program loaded into it.

→ 2.4.2 Process Termination

Q. Explain process termination. **J**

- When process finishes the execution of last statement, it terminates. After this it requests the operating system to delete it by using the `exit()` system call. Just then, the process may return an integer to its parent process by using `wait()` system call.
- Then all allocated resources to the process like physical and virtual memory, open files, and VOs buffers are freed by the operating system. Termination can happen in other situations also.
- By executing the suitable system call, any process can terminate the other process. Normally, parent of the process to be terminated, executes this system call.

Process Management

- Once Twister users could randomly kill other's jobs. It is necessary that a parent requires knowing the identity of its children.
- Thus, when a process creates a new process, the identity of the child process is passed to the parent.
- Any of the child processes execution can be terminated by child for a variety of reasons, such as these:
 - If the allocated resource usage is exceeded by child. (The parent must have a means to examine the state of its children.)
 - There is no longer requirement of allocated task to the child.
 - In many systems, If the parent is terminated then OS does not permit a child to carry on execution.
- For example, VMS system does not permit a child to carry on execution after its parent process is terminated.
- If all children of particular parent are terminated then it is called as cascading termination. After normal or abnormal termination of parent, cascaded termination is carried out in much system. It is initiated by the OS.
- In UNIX, `exit()` system call is used to terminate the process. The `wait()` system call is used by parent to wait until child is terminated.
- The terminated child's process identifier is returned by the `wait()` system call. Because of this process identifier, parent process comes to know about which child is terminated.
- After the termination of the parent, the `rmif` process is assigned as parent to all terminated children.
- This `init` process collects all terminated children's status and execution statistics after termination of their parent.

2.5 Process Control Block

→ (Dec. 14)

Q. Explain role of process control block.

MU - Dec. 2014. 5 Marks

Q. Explain process control block. **J**

- Any process is identified by its process control block (PCB). PCB is the data structure used by the operating system to keep track on the processes.
- All the information associated with process is kept in process control block. There is separate PCB for each process.

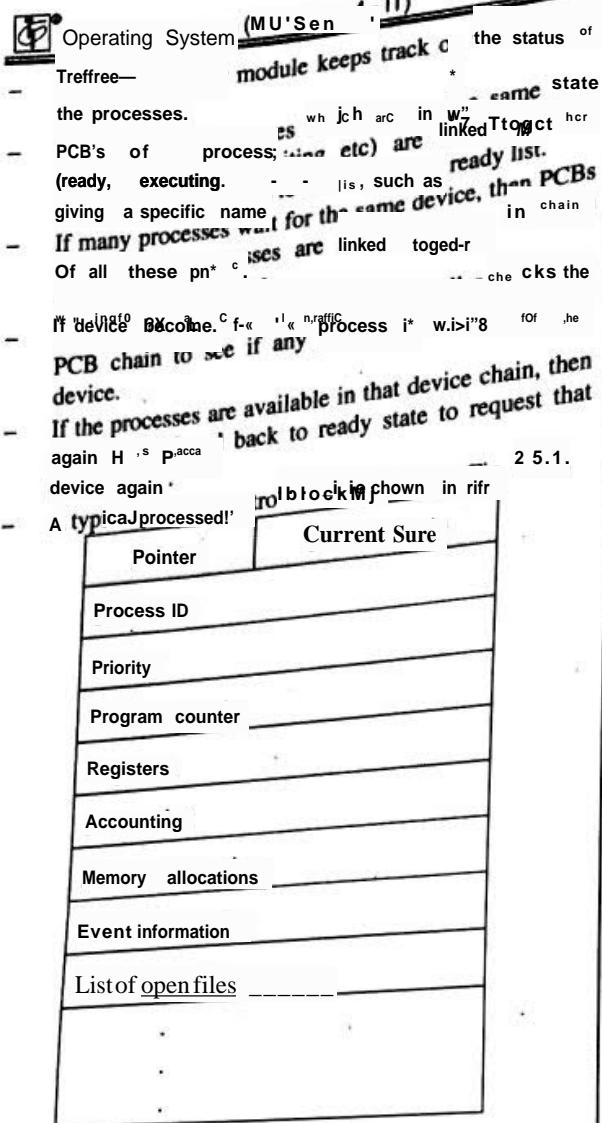


Fig. 15.1 : Process control block

»■ Pointer

This field points to other process's PCB. The scheduling list is maintained by pointer.

» Current state

Currently process can be in any of the state from new, ready, executing, waiting etc as described above.

» Process ID

Identification number of the process. Operating system assign this number to process to distinguish it from other processes.

» Priority

Different process can have different priority. Priority field indicate the priority of the process.

ram counter

PrOS CPU is given to other processes
 After c o nic X' s*, ch previously executing process
 When turn of 0.6 P dlocated to it. Program counter
 address of the context switch.

j5. <— **

Refill

It includes general purpose register, index registers, stack pointers and multipliers etc. number of registers and type of register differs as per the architecture of computer. The state information stored after interrupt is again used by the process to resume the execution.

» Account

nt fl calculating the process's priority

Information

relative to other Processes.

This may include accounting information about CPU members. It also includes amount of time used, time limits, process and so on.

Memory allocation

This information may include the value of base and limit.

It register. It includes paging, segmentation information depending on the memory system Address space allocated to the process etc.

» Event information

For a process in the blocked state this field contains information concerning the event for which the process is waiting.

» List of open files

Files opened by the process.

After creation of the process, hardware registers and flags are set as per the details supplied by loaders or linker. At any time the process is blocked, the processor registers content are generally placed on the stack and the pointer to the respective stack frame is stored in the PCB. In this fashion, the hardware state can be restored when the process is scheduled and resume execution again.

2.6 Process States and Process State Transition Diagram

Q. Draw and explain process state transition diagram. During execution, process changes its state. Pro®⁹

state contains five states: current activity of the process. Pro®⁹

- Each process remains in one of these five states. There is a queue associated with each state of the process.
- Process resides on that queue as per the state in which it resides. The states are listed in Fig. C2.3.

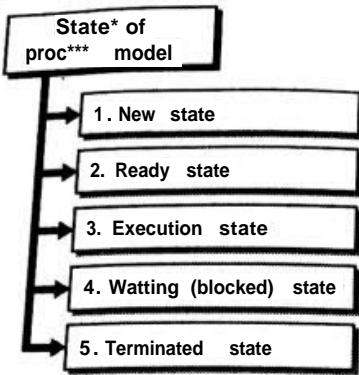


Fig. C23 : States of Process Model

-> 1. New state

The new process being created.

-> 2. Ready state

A process is ready to run but it is waiting for CPU being assigned to it.

*♦ 3. Executing state

A process is said to be in running state if currently CPU is allocated to it and it is executing.

■ 4. Waiting (blocked) state

A process can't continue the execution because it is waiting for event to happen such as I/O completion. Process is able to run when some external event happens.

-> 5. Terminated state

The process has completed execution.

The process state transition diagram is shown in Fig. 2.6.1.

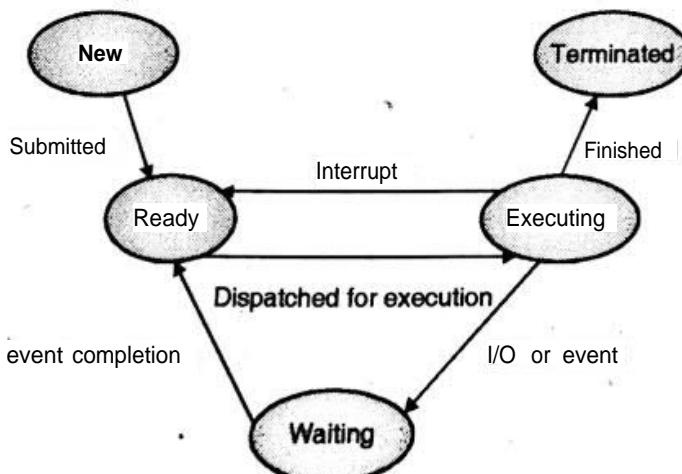


Fig. 2.6.1 : Process state transition diagram

When the process is created, it remains in new state. After the process admitted for execution, it goes in ready state.

A process in this state, wait in the ready queue. Scheduler dispatches the ready process for execution i.e. CPU is now allocated to the process.

When CPU is executing the process, it is in executing state. After context switch, process goes from executing to ready state.

If executing process initiates an I/O operation before its allotted time expires, the executing process voluntarily give up the CPU.

In this case process transit from executing to waiting state. When the external event for which a process was waiting happens, process transit from waiting to ready state.

When process finishes the execution, it transit to terminated state.

2.7 Process vs. Thread

- A thread is a single sequence stream within a process. Threads are also called as lightweight processes as it possess some of the properties of processes. Each thread belongs to exactly one process.
- In operating system that support multithreading, process can consist of many threads.
- These threads run in parallel improving the application performance. Each such thread has its own CPU state and stack, but they share the address space of the process and the environment.
- Threads can share common data so they do not need to use interprocess communication.
- Like the processes, threads also have states like ready, executing, blocked etc. priority can be assigned to the threads just like process and highest priority thread is scheduled first.
- Each thread has its own Thread Control Block (TCB). Like process context switch occurs for the thread and register contents are saved in TCB.
- As threads share the same address space and resources, synchronization is also required for the various activities of the thread.

o- Difference between **thread and process**

Parameters		process	Thread
Sr. No*		Program in execution called as process. heavy weight process.	Thread is P ¹ of process, hence so called light weight process.
1. Definition		Process context switch time as compared to thread context switch.	Thread context switch takes less time as compared to process context switch.
2. Context switch		Process creation takes more time as compared to new thread creation.	Thread creation takes less time compared to new process creation.
3. Creation		New Process creation takes more time as compared to new thread creation.	New thread creation takes less time compared to new process creation.
4. Termination		New Process termination takes more time as compared to new thread termination.	New thread termination takes less time compared to new process termination.
5. Execution		Each process executes the same code but has its own memory and file resources.	All threads can share same set of open files, child processes.
6. Implementation		If implementation is process based, then blocking of one process cause the blocking of other server process until the first process unblocked and these are not allowed to execute until blocked process is unblocked.	In multithreaded server implementation, if one thread is blocked and waiting, second thread in the same process could execute.
7. Multiple resources		Multiple redundant processes use more resources than multiple threaded process.	With compare to multiple redundant process, multiple threaded processes need fewer resources.
8. Address space		Context switch flushes the MMU (TLB) registers as address space of process changes.	No need to flush TLB as address space remains same after context switch because threads belong to the same process.

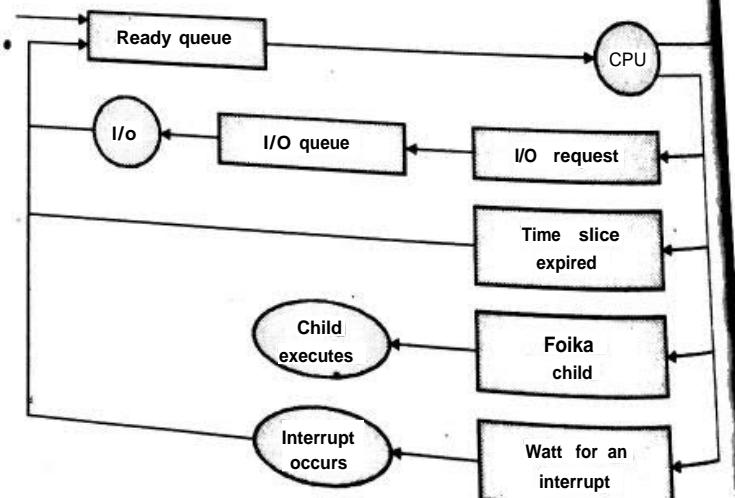
Syllabus Topic : Process Scheduling

2.8 Process Scheduling

In order to improve CPU utilization, multiprogramming concept makes available several processes ready for execution.

2.8.1 Scheduling Queues and Schedulers

- System places all the ~~int~~ **ready processes** in the **job queue**.
- When ~~are~~ **idle**, **side** **job queue** for CPU, they kept in **util** **and wait** **queue**.
- AW from the job **queue**.
- operating system also has** **processes** **device, the** **device, the** **respecti** **vedevice** *** Pedlar**



Fi

8 2 8 1 1 Of n Uin8 diaBran , representation
process scheduling

After allocation of CPU to the process start the execution. While executing the process, one of the numbers of events could occur.

Until the process finishes the execution, it travels between various scheduling queues. The operating system selects the processes from queues based on some policy.

This selection is carried out by the program called as scheduler. There are three types of schedulers to schedule a process. These are shown in Fig. C2.4.

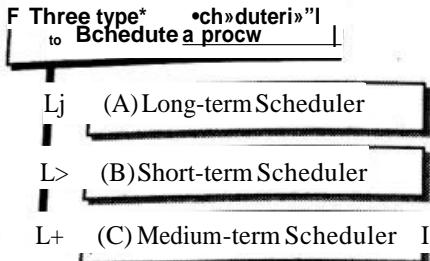


Fig. C2.4 : Types of scheduler

-> 2.8.1(A) Long-term Scheduler

Q. Explain long-term scheduler.

- When programs are submitted to the system for the purpose of processing, long term scheduler comes to know about it.
- Then its job is to choose processes from the queue and place them into main memory for execution purpose.
- CPU bound processes require more CPU time and less I/O time until execution completes.
- On the contrary, I/O bound processes use up more time in doing I/O and require less CPU time for computation.
- The main job of the long term scheduler is to provide a balance mix of I/O bound and CPU bound jobs.
- The number of processes in memory for execution and degree of multiprogramming is related to each other.
- More number of processes in memory for execution indicates degree of multiprogramming is high. Long term scheduler controls the degree of multiprogramming.
- If the average rate of new process creation and average departure rate of processes leaving the system is equal then degree of multiprogramming is steady.
- The Long term scheduler is not present in timesharing operating systems.

When process state transition take place from new to ready, then there long term scheduler come into picture for scheduling purpose.

→ 2.8.1(B) Short-term Scheduler

Q. Explain short-term scheduler.

- Processes which are in ready queue wait for CPU. A short term scheduler chooses the process from ready queue and assigns it to the CPU based on some policy.
- These policies can be First Come First Served (FCFS), Shortest Job First (SJF), priority based and round robin etc, Main objective is increasing system performance by keeping the CPU busy.
- It is the transition of the process from ready state to running state. Actual allocation of process to CPU is done by dispatcher.
- Short term scheduler is faster than long term scheduler and should be invoked more frequently compare to long term scheduler.

■ 2.8.1(C) Medium-term Scheduler

Q. Explain medium-term scheduler.

- If the degree of the multiprogramming increases, medium-term scheduler swap out the processes from main memory.
- The swapped out processes again swapped in by medium-term scheduler.
- This is done to control the degree of multiprogramming or to free up a memory.
- This is also helpful to balance the mix of different processes, some time sharing operating system have this additional scheduler.

2.8.1(D) Comparison of Three Schedulers

Q. Compare the functions of different types of schedulers.

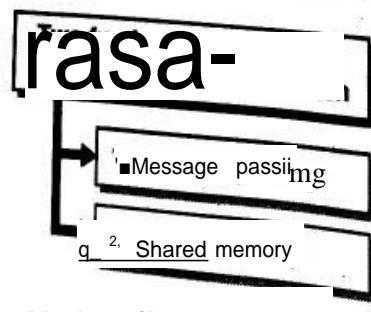
Sr. No.	Long-term Scheduler	Short-term Scheduler	Medium-term Scheduler
1,	Selects processes from the queue and loads them into memory for execution.	Chooses the process from ready queue and assigns it to the CPU.	Swap in and out the processes from memory.

Sr. No.	Long-term Scheduler	Short-term Scheduler	Medium-term Scheduler
2.	Speed is less than the short term scheduler.	Speed is very fast and invoiced frequent ly than long scheduler.	Speed is in between both short term scheduler and long scheduler.
3.	Transition of process state from ^{Ne*} to Ready.	Transition of process state from Ready to Executing.	No process state transition'
4.	Not present in time sharing system.	Minimal in time sharing system.	Present in time sharing system.
5.	Supply a reasonable mix of jobs, such as I/O bound and CPU bound	Select a new process to allocate to CPU frequently.	Processes are swapped in and out for balanced process mix.
6.	It controls degree of multiprogramming through placing processes in ready queue.	It has control over degree of multiprogramming as it allocates processes to CPU for execution.	Reduce the degree of multiprogramming by swapping the processes in and out.
7.	It is also called as job scheduler.	It is also called as CPU scheduler.	It is used for swapping.

Syllabus Topic: interprocess Communication

2.9. J) ~MssComm

unicat



W> Cts8

Modus of Interpret

z2

Om "uni cati , ri

action and communication requirements should be satisfied with each other. Synchronization is required to achieve the mutual exclusion.

Independent processes do communicate with each other but cooperate. Cooperative processes communicate through shared memory or passing.

Message naming provides both functions, using passing has the further benefit that it lends itself to implementation in distributed systems as well as multiprocessor and uniprocessor shared-memory systems.

Following are the two primitives used in passing:

1. send (destination, message)
2. receive (source, message)

This is the minimum two operations required by processes to send and receive the messages. A process sends data in the form of a message to another process indicated by a destination. A process receives data by executing the receive primitive, indicating the source and the message.

Communication by sending and receiving messages require synchronization. The receiver cannot receive a message until it has been sent by another process.

The sending process is blocked until the message is received, or it is not after the send primitive is issued by process. Similarly, when a process issues a receive primitive, there are two possibilities:

- o Previously sent message is received and process continues.
- o there is no waiting message, then either process is blocked until a message arrived or the process continues to execute, abandoning the attempt to receive.

Thus

* both the sender and receiver can be blocked.

Then

Particularly in synchronous communication

«@*will be implemented: y although

- o **The blocking send and blocking receive :** Until the message is handed over, the sender and receiver both gets blocked.
- o **The nonblocking send and blocking receive :** After sending the message, the sender continues its work but receiver remains blocked until message is arrived to it. This permits a process to send one or more messages to a multiple destinations as quickly as possible. Here receiver is in need of message so that it can resume the execution. So it gets blocked until message arrives.
- o **Nonblocking send, nonblocking receive :** Neither party is required to wait.
- o **The nonblocking send and non blocking receive:** Both sender and receiver will not wait and both will continue the work.
- Message passing system should give guarantee that messages will be correctly received by receiver.
- Receiver sends acknowledgement to sender after receiving the message.
- If acknowledgement not received in defined time then sender resend the message. It also offers authentication service.

→ 2.9.2 Shared Memory

- Cooperating processes require an Interprocess Communication (IPC) mechanism that will allow them to exchange data and information.
- In the shared-memory model, a region of memory that is shared by cooperating processes is established.
- Processes can then exchange information by reading and writing data to the shared region.
- In the message passing model, communication takes place by means of messages exchanged between the cooperating processes.

Syllabus Topic : Multithreading

2.10 Multithreading

- A thread is a single sequence stream within a process. Threads are also called as lightweight processes as it possess some of the properties of processes. Each thread belongs to exactly one process.

- In operating system that support multithreading, process can consist of many threads.
- These threads run in parallel improving the application performance. Each such thread has its own CPU state and stack, but they share the address space of the process and the environment.
- Threads can share common data so they do not need to use interprocess communication. Like the processes, threads also have states like ready, executing, blocked etc. priority can be assigned to the threads just like process and highest priority thread is scheduled first.
- Each thread has its own Thread Control Block (TCB). Like process context switch occurs for the thread and register contents are saved in TCB.
- As threads share the same address space and resources, synchronization is also required for the various activities of the thread.

2.11 Types of Threads

- Q. Explain user level and kernel level threads with advantages and disadvantages.

Threads can be implemented in two ways.

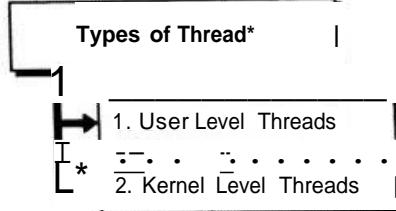


Fig. C2.6 : Types of Threads

→ 1. User Level Threads

- In user level implementation, kernel is unaware of the thread. In this case, thread package entirely put in user space. Java language supports threading package.
- User can implement the multithreaded application in java language.
- Kernel treats this application as a single threaded application. In a user level implementation, all of the work of thread management is done by the thread package.
- Thread management includes creation and termination of thread, messages and data passing between the threads, scheduling thread for execution, thread synchronization and after context switch saving and restoring thread context etc.

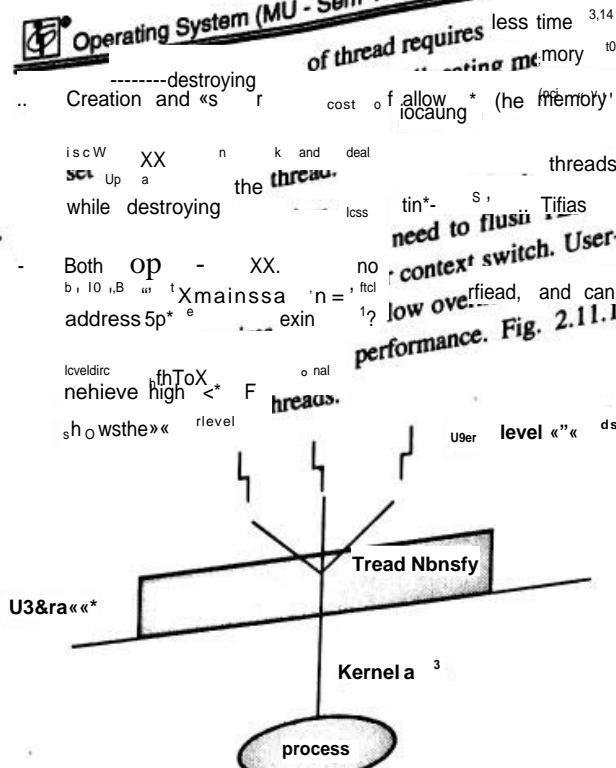


Fig. 2.11.1 : User level threads

Advantages of user level threads

- Thread switching does not require flushing of CPU register doing CPU accounting. Only need to be stored and reloaded again.
- user level threads are platform independent, they can execute on any operating system.
- Scheduling can be as per need of application.
- Thread management are at user level and do not require thread library. So kernels burden is taken by threading package. Kernels time is saved for other activities

Disadvantages of user level threads

- If one thread is blocked on I/O, entire process gets blocked.
- The applications where after blocking of one thread other requires to run in parallel, user level threads are of no use.

→ 2 Kernel Level Threads

- In this, threads are implemented in operating system's kernel. The thread management is carried out by kernel.
- All these thread management activities are carried out in kernel space. So thread context and process context switching becomes same.

Application can be written as multithreaded. Application are supported as follows:

processes are generated by TT more S

Kernel threads are more than the user threads.

Kernel threads are generated by kernel.

Address schedule another thread of

The kernel tough one thread in a process is blocked even if it does not block the thread.

Block of

Kernel can simultaneously schedule multiple threads from the same process.

Kernel level threads

Kernel intervention.

Contest switch are generated requires more time.

Kernel threads create and manage than the user threads.

Process - Multithreading Models

Syllabus Top

2.1.2 Multithreading Models

models.

Q. Explain van

3 * 03

- Herings the advantages of user level and fe, X threads, a hybrid threading model using of threads can be implemented.
- The Solaris operating system supports this hybrid model.
 - In this implementation, all the thread management functions are carried out by user level thread package user space. So operations on thread do not require kernel intervention.
 - The advantage of hybrid model of the threads is that the applications are multithreaded then it can take advantage of multiple CPUs if they are available. 0 threads can continue to make progress even if kernel blocks one thread in a system function,
 - There are three types of multithreading models. These are shown in Fig. C2.7.

Three types of multithreading models

1, One to one

2. Many to one

3, Many to many

Fig. C2.7 : Types of multithreading models

→ 1. One to one model

- In this model relationship between user level thread and kernel level thread is one to one.
- It means that there is mapping of a single user-level thread to a single kernel-level thread.
- Because of such type of relationship multiple threads executes in parallel leading to more concurrency.
- However, since it is needed to create kernel thread for every new creation of user thread.
- So application performance will be degraded. Windows series and Linux operating systems try to minimize this problem by restricting the expansion of the thread count. OS72, Windows NT and windows 2000 use one to one relationship model,

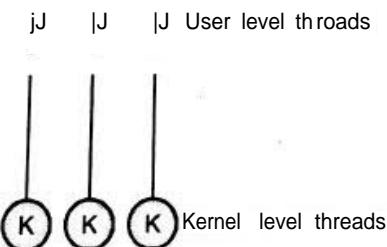


Fig. 2.12.1

→ 2. Many to one model

- In this model relationship between user level thread and kernel level thread is many to one.
- It means that there is mapping of many user-level threads to a single kernel-level thread.
- In this model management is done in user space. When one thread makes a system call for blocking, the entire process gets blocked.

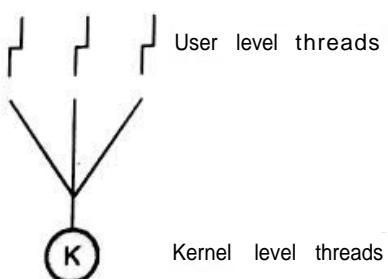


Fig. 2.12.2 : Many to one model

- At a time only one thread can access the Kernel thread at a time, so many other threads cannot execute in parallel on multiple processors.

- Therefore concurrent execution of threads cannot be achieved.

- This type of relationship facilitates an effective context-switching environment, easily implementable even on simple kernels with no thread support

→ 3- Many to Many Model

- Many to many association exist between user level thread and kernel level thread in this model. It means that more number of user-level threads are allied to equal or less number of kernel-level threads.
- The necessity of altering code in both kernel and user spaces leads to a level of complexity not present in the one to one and many to one model.
- Like many-to-one model, this model offers an efficient context-switching environment as it repels from system calls.
- The keen complexity offers the potential for priority inversion and suboptimal scheduling with minimum coordination between the user and kernel schedulers.

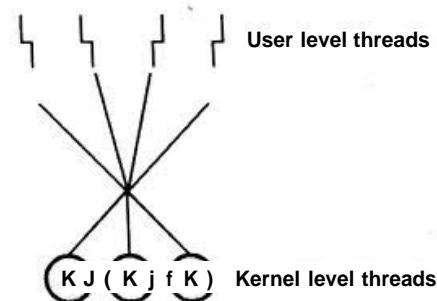


Fig. 2.12.3 : Many to Many model

Syllabus Topic : Thread Librads*

2.13 Thread Libraries

q Write note on thread libraries.

- Programmers can create and manage the threads using API that are provided by thread library.
- Thread library can be implemented at user level where it runs in user space without kernel support. Hence calling a function in library is local function call in user space.
- Kernel level library is supported by operating system. The data structures and code remains in kernel space.

- Hence any invocation of function in library » **call to the kernel but no** bread libraries that are in use. Following are the main libraries today.

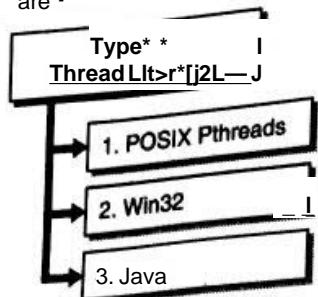


Fig. C2.8 : Thread libraries

» 2.13.1 POSIX Pthreads

- **Pthreads** is the threads extension of the POSIX standard that can be provided as either a user-level or kernel-level library.
- Pthreads define API for the creation of threads and synchronization. As it is a specification, the design of operating can choose their own way of implementation of specification.
- Solaris, Linux, Mac OS X, and Tru64 UNIX are the systems that implements Pthreads specification. For all the types of Windows operating systems, a shareware implementation is available in public domain.

» 2.13.2 Win32 Threads

- It is kernel-level library available on windows system. The approach for creation of threads using Win32 is similar to Pthreads.
- Just as Pthreads, in Win32 CreateThread() function is used to create the threads.
- The required set of attributes for the thread is passed to this function. Security information, stack size, and a **flag which is set to specify if the thread is to start in a suspended state** are the examples of these attributes.

» 2.13.3 Java Threads

- Java threads can be directly created out in Java programs. The creation and management of threads can be carried out in Java programs.
- **As most of the time, JVM is runs on the top of a host operating system** the threads are created on the host system.

Windows systems. Java threads are

implemented using the Win32 and *0 in UNIX and Linux systems.

The API to manage and create the threads. The `new Thread()` method in java program is the class `Thread`.

least single thread is supported.

In first approach of thread creation, a new class is created which is the `Thread`.

Thread Class and then overrides its `run()` method. In second approach a class is defined that implements the `Runnable` interface.

Syllabus Topic : Threading

2.14 Threading Issues

- Q. What are different threading issues?

Following are the threading issues that needs to be considered with multithreaded programs.

Threading Issues

1. The fork() and exec() System Calls
2. Cancellation
3. Signal Handling
4. Thread Pools
5. Thread-Specific Data
6. Scheduler Activations

Fig. C2.9 : Threading issues

» 2.14.1 The fork() and exec() System Calls

- A separate duplicate process is created with **fork system call**. The meaning of `fork()` and `exec()` system calls is different in multithreaded program.

in UNIX Versions when single thread of program () then it duplicate all the threads of program.

LtLT." " UN". -r »-pi—«*

- When thread invokes exec() system call, the whole process including all the threads will be replaced by program that is given as parameter to exec() call.
- The use of particular version of forkf) is carried out as per application need.
- If it is necessary to call exec() just after forking then no need to duplicate all the threads. If exec() is not called just after forking then duplication of all the threads is required.

→ 2.14.2 Cancellation

- Termination of thread before it completes the execution is thread cancellation. If simultaneously many threads are searching through database and one thread return the result then cancellation of remaining threads is done.
- The thread chosen for cancellation is called as target thread. Following two scenarios are there to cancel the target thread.
 - o **Asynchronous cancellation.** Instant termination of the target thread is done by one thread
 - o **Deferred cancellation.** The periodic check is carried out by target thread to decide whether it should terminate, permitting it an opportunity **to terminate itself in an orderly fashion.**
- The asynchronous cancellation is troublesome in situation if cancellation of thread is carried out to which system resources have been allocated.
- Asynchronous cancellation is also difficult in case thread is in middle stage of updating data that it is **shared with other threads.** Although operating system will reclaim the resources from cancelled thread some resources will be there with cancelled thread and cannot be reclaimed.
- In deferred cancellation, thread to be cancelled checks the flag to know whether it should be cancelled or not. **This always allows safe cancellation.**

→ 2.14.3 Signal Handling

In UNIX, process is **informed by signal when some event is occurred.** The signals can be synchronous or asynchronous depending on source and reason for event being signal. Following pattern is followed by synchronous or asynchronous signals.

- When particular event occurs, the signal is generated.
- This generated signal is **undelivered** to the process.

- The delivered signal should be handled

Following are the examples of synchronous signals.

- Illegal memory access.
- Division by zero.

If the executing program carries out above two actions then generated signals are delivered to the same process which performed above actions. Hence, these signals are called as synchronous signals. When signals are due to events of external process then the executing process receive asynchronous signals

Following are examples of asynchronous signals.

- Terminating process by applying some specific keystrokes.
- Timer expires.

Both types of signals are handled either by default signal handler or user defined handler. If program is multithreaded program then following options are available for delivery of these signals to the process..

- Signal can be delivered to the thread to which it is applicable.
- Signal can be delivered to every thread in the process.
- Signal can be delivered to specific threads in the process.
- Allocate any one particular thread to take delivery of all signals for the process

The way signals are delivered depends on types of signals generated.

→ 2.14.4 Thread Pools

- There should be bound on number of threads active in system. If unlimited number of threads exists in system then many resources will be consumed by them.
- Thread pool is solution on this issue. When process starts up then some number of threads are created. These threads are then placed in pool and they waits or work.
- If server comprising these threads receives request, it **passes it to any of the available waiting thread in pool.**
- Once thread completes the requested work, it returns to pool and waits for other work.
- **If thread is not available in pool then server has to wait till] thread will be available.**

Operating System (MU - Sem 4 - IT)

- Following benefits are provided due to thread pool:
- Instead of waiting to ^{rent}³ thread servicing a request with an existing thread is always faster.
 - A thread pool provides a way to manage threads based on number of threads that exist at any point of time. It helps in reducing the number of parallel threads.

Data

» 2.14.5 Thread-Specific

- Always threads of a process share the data that is available to all threads. Surely, this is not suitable for multithreaded programming.
- On the other hand, if threads in a process require their own copy of data, data is a thread-specific data.

→ 2.14.6 Scheduler Activations

- It is important to consider common coordination between the kernel and the thread library, which may be needed by the kernel level threads.
- This coordination permits the number of kernel threads to be dynamically adjusted to guarantee the best performance.
- Several systems implementing either the many-to-many or user and kernel level model of threads place an intermediate data structure between the user and kernel threads.
- This data structure is called as lightweight process, or LWP. For user level library LWP is virtual processor to schedule the application to run.
- LWP is attached with kernel thread and OS schedules LWPs as kernel threads to run on physical processor. With blocking of kernel thread, LWP also blocks and vice versa. User thread associated with LWP also blocks.
- Scheduler activation is the method of communication between the user-thread library and the kernel.
- The kernel offers an application with some number of virtual processors (LWPs). It can schedule user threads onto available virtual processor.

2.15 Process Scheduling

2.15.1 Scheduling Decisions

- Scheduling decisions are made in following situations.
- After creation of new process. Decision is needed for parent or child process for execution.
 - When any process is selected.
 - When process is blocked due to I/O or because of other reason.
 - On the occurrence of interrupt.

M Type scheduling

What is difference between preemptive and non-preemptive scheduling?

Types of scheduling algorithms

1. Non-Preemptive
2. Preemptive

Fig. C2.10 : Types of scheduling algorithms

1. Non-Preemptive

Non-preemptive algorithms are designed so that once a process is allocated to CPU, it does not free CPU until it completes its execution.

2. Preemptive

Preemptive algorithms allow taking away CPU from a process during execution. If highest priority process arrives in the system, CPU from currently executing low priority process is allocated to it. It ensures that always highest priority process will be executing.

Dispatcher and dispatch latency

- Short term scheduler allocates CPU to ready processes based on some criterion. I.e., the traits of the processes from ready queue, running queue, or waiting queue.

- Actual allocation of process to CPU is done by dispatcher. Short term scheduler is faster than long term scheduler and should be invoked more frequently compare to long term scheduler.
- **The amount of time dispatcher needed to bring to stop** one process and starts running of other process is called as dispatch latency.

 Syllabus Topic : Scheduling Criteria

2.16 Scheduling

2.16.1 Scheduling Criteria

Q. What are the criteria for evaluation of scheduling algorithm performance?

- In multiprogramming, many programs remain in memory at the same time.
- Processes carry out I/O operations while performing the execution. Since I/O operations generally consume more time to accomplish with compare to CPU instructions, multiprogramming systems hand over the CPU to another ready process whenever any process invokes an I/O operation.
- Short term scheduler allocates CPU as per some policy called as scheduling algorithms.
- The main aim of the scheduling is to improve performance by keeping CPU busy all the time.
- Criteria for performance evaluation of the scheduling strategy is :
 - o **CPU Utilization** : It is amount of time CPU remains busy.
 - o **Throughput** ; Number of jobs processed per unit time.
 - o **Turnaround time** : Time elapsed between submission of job and completion of its execution.
 - o **Waiting time** : Processes waits in ready queue to get CPU. Sum of times spent in ready queue is waiting time.
 - o **Response Time** : Time from submission till the first response is produced.
 - o **Fairness** : Every process should get fair share of CPU time.

 Process Management

 Syllabus Topic: Scheduling Algorithms

2-1 S-2 Scheduling Algorithms

Scheduling algorithm.

- (A) First in First out (FIFO)
- (B) Shortest Job First (SJF)
- (C) Priority Scheduling
- (D) Round Robin Scheduling
- (E) Multilevel Queue Scheduling
- (F) Multilevel Feedback-Queue Scheduling

Fig. C2J 1 : Scheduling algorithms

→ 2.16.2(A) First in First Out (FIFO)

Q. Explain FIFO scheduling algorithm.

- This is a Non-preemptive scheduling algorithm, FIFO strategy allocates the CPU to processes in the order of their arrival.
- This algorithm treats ready queue as FIFO. A process does not give up,
- CPU until it either terminates or performs I/O, If the longer job assigned to CPU then many shorter jobs has to wait. As long processes can hold the CPU, this algorithm gives fewer throughputs.
- This algorithm can be easily implemented using FIFO queue. In time sharing system every user should get the CPU time at regular interval. So FIFO algorithm is not useful in such situations.
- FIFO algorithm is inappropriate for interactive systems; large fluctuations in average turnaround time are possible.
- If we assume that arrival time is zero. Consider the following example :

Process	Burst time
P1	24
P2	03
P3	12

Assume that processes arrive in the order P1, P2, and P3. The Gantt chart shows the result.

	P1	P2	P3	32
0	24	2?		

$$\text{Turnaround time for } P2 = (3 - 0) = 3$$

$$\text{Turnaround time for } P3 = (8 - 0) = 8$$

$$\text{Turnaround time for } P1 = (32 - 0) = 32$$

$$\text{Turnaround time for } P1 = (3 + 8 + 32) / 3 = 14.33$$

Turnaround time

Average Jum

Waiting time for P2

0

= 3

Waiting time for P2

= 8

Waiting time for P3

= 8

Waiting time for P3

= $(0 + 3 + 8) / 3 = 4.66$

Average Waiting time

Turnaround time and average waiting time varies with order of arrival.

2.16.2(B) Shortest Job First

Q. Explain SJF and SRTN scheduling algorithms.

- Shortest Job First (SJF) may be preemptive or non preemptive.

the jobs so as to run the shortest jobs first e response time.

(SJF) improves the average

ned in order of increasing job

Ready queue is maintained in the ready

queue based on execution time.

queue based on execution time. SJF minimizes the average

wait time because it gives service to less execution time

processes before it gives service to large execution time

processes.

- While it minimizes average wait time, it may punish processes with high execution time.

- If shorter execution time processes are in ready list, then processes with large service times tend to be left in the ready list while small processes receive service.

- It may happen in extreme case that always short execution time processes will be served and large execution time processes will wait indefinitely.

- This starvation of longer execution time processes is the limitation of this algorithm.

- Consider the example discussed for FCFS algorithm. In this example it is assumed that SJF is non preemptive.

- If the order of arrival is P2, P3, P1, the order of execution time will be 3, 5, and 24. There is significant reduction in average turnaround time and average waiting time.

Consider the following example of SJF algorithm.

Process	Arrival time	Burst time
P1	0	10
P2	1	5
P3	2	8
P4	3	15

P1	P2	P3	P4
10	5	23	38

$$\text{Turnaround time for } P1 = (10 - 0) = 10$$

$$= (15 - 1) = 14$$

$$= (23 - 2) = 21$$

$$= (38 - 3) = 35$$

$$= (10 + 14 + 21 + 35) / 4 = 21$$

Average Turnaround time

for P1

Waiting time for P1

0

= (10 - 1) = 9

Waiting time for P2

= (15 - 2) = 13

Waiting time for P3

= (23 - 3) = 20

Waiting time for P4

= (38 - 4) = 34

Average waiting time

In Preemptive SJF, if newly arrived process execution time with compare

current executing process, then the CPU will be given to newly arrived process.

The preemptive Shortest Job First (SJF) is called

Shortest Remaining Time Next Scheduling (SRTN)

Consider the above example for preemptive SJF.

Job First (SJF). The Gantt chart shows the result.

P1	P2	P3	P1	P4
0	1	6	14	23

- P1 arrives at time 0, so get the CPU as it is only process in the queue. At time 1, process P2 arrives having execution time 5 (less than remaining time 9 of P1). P1 is preempted and P2 is scheduled. At time 2, process P3 arrives having execution time 8 which is greater than remaining time of P2.
- So P2 will complete the execution. After P2, process P3 will execute as its execution time is less than P1 (9) and P4 (15). Process P1 has remaining execution time 9. So it is scheduled next. Finally P4 will complete the execution.

$$\text{Turnaround time for } P1 = (23 - 0) = 23$$

Turnaround time for P3 = $(14 - 2) = 12$
 Turnaround time for P4 = $(38 - 3) = 35$
 Average Turnaround time = $(23 + 5 + 12 + 35) / 4 = 18.75$
 Waiting time for P1 = $(14 - 1) = 13$
 Waiting time for P2 = $(1 - 1) = 0$
 Waiting time for P3 = $(6 - 2) = 4$
 Waiting time for P4 = $(23 - 3) = 20$
 Average Waiting time = $(13 + 0 + 4 + 20) / 4 = 18.75$

→ 2.16.2(C) Priority Scheduling

Q. Describe priority scheduling algorithm.

- In priority scheduling, each process has a priority which is a integer value assigned to it.
- Smallest integer is considered as highest priority and largest integer is considered as lowest priority. Always highest priority process gets the CPU.
- In some system, largest number is treated as a highest priority and it depends on implementation.
- If priorities are internally defined then some measurable quantity such as time limits, memory requirements, the number of open files, and the ratio of average I/O burst to average CPU burst are used to compute the priorities.
- External priorities are assigned on the basis of factors such as importance of the process, the type and amount of funds being paid for computer use, the department sponsoring the work, and other, often political, factors. All these factors are not related to the operating system.
- Preemptive and non preemptive SJF is a priority scheduling where priority is the shortest execution time of job.
- In this algorithm, low priority processes may never execute. This is called **starvation**.

Solution to this starvation problem is **aging**. In aging as time progresses, increase the priority of the process so that lowest priority processes gets converted to highest priority gradually.

Consider the following example for priority scheduling.

The Gantt chart shows the result.

Process	Burst time	Priority
1 - P1	12	4
P2	10	5
P3	5	2
P4	4	1
P5	3	3

P4	P3	P5	P1	P2
0	4	9	12	24

Turnaround time for P1 = $(24 - 0) = 24$
 Turnaround time for P2 = $(34 - 0) = 34$
 Turnaround time for P3 = $(9 - 0) = 9$
 Turnaround time for P4 = $(4 - 0) = 4$
 Turnaround time for P5 = $(12 - 0) = 12$
 Average Turnaround time = $(24 + 34 + 9 + 4 + 12) / 5$
 $= 16.6$

Waiting time for P1 = 12
 Waiting time for P2 = 24
 Waiting time for P3 = 4
 Waiting time for P4 = 0
 Waiting time for P5 = 9

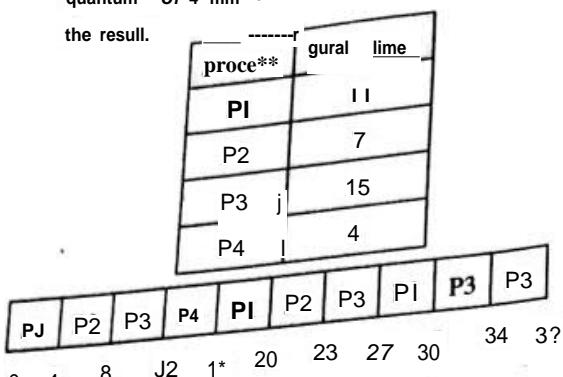
Average Waiting time = $(12 + 24 + 4 + 0 + 9) / 5 = 9.8$

→ 2.16.2(D) Round Robin Scheduling

Q. With the help of example, explain Round Robin Scheduling algorithm.

- Round Robin Scheduling is designed especially for time-sharing systems where many processes get CPU on time sharing basis.
- In this algorithm, a small unit of time called time quantum is defined.
- CPU is allocated to each process for this time quantum period of time. To implement this scheduling, ready queue treated as FIFO queue.
- The new processes go to the tail of queue and each time CPU chooses the process from head of queue.
- When time quantum expires, context switch occurs and CPU switches to other process which is scheduled next.
- The time quantum is fixed and then processes are scheduled such that no process get CPU time more than one time quantum.
- If process is executing and request for I/O the process goes in waiting (blocked) state.
- After the completion of I/O, process again gets added at the tail of ready queue. The time quantum should not be very small or very large.
- If the time quantum is very large, the algorithm will behave just like FCFS.
- A smaller time quantum increases context switches leading to performance degradation (less throughput) as context switches elapses more time.

- Consider the following example. If we use quantum of 4 ms then the Gantt chart shows the result.



$$\begin{aligned}
 \text{Turnaround time for P1} &= (30 - 0) = 30 \\
 &= (23 - 0) = 23 \\
 \text{Turnaround time for P2} &= (37 - 0) = 37 \\
 \text{Turnaround time for P3} &= (16 - 0) = 16 \\
 \text{Turnaround time for P4} &= (23 + 37 + 16) / 4 = 26.5
 \end{aligned}$$

$$\text{Average Turnaround time} = (23 + 26.5) / 4 = 19.875$$

$$\text{Waiting time for P1} = 0 + (16 - 4) + (27 - 20) = 19$$

$$\text{Waiting time for P2} = 4 + (20 - 8) = 16$$

$$\text{Waiting time for P3} = 8 + (23 - 12) * (30 - 27) = 22$$

$$\text{Waiting time for P4} = 12$$

14.25

$$\text{Average Waiting time} = (19 + 4 + 22 + 12) / 4 = 16.25$$

→ 2.16.2(E) Multilevel Queue Scheduling

Q. Explain multilevel queue scheduling algorithm.

- Sometimes it is necessary to categorize the processes into different groups. For example, separation is made between interactive processes and batch processes.
- The response-time need of these two processes can be dissimilar. So these processes can have different scheduling requirement. Also, interactive processes may have higher priority over batch processes.
- In multilevel queue scheduling algorithm, there are multiple ready queues. The allocation of process to the particular queue depends on property of that process such as memory size, priority of the process, or its type.
- Scheduling algorithm for different queues can be different. The round robin scheduling algorithm can be used for interactive processes queue and first come first serve scheduling algorithm can be used for batch processes queue.
- In addition to scheduling for each queue, there must be scheduling between the different queues. The fixed-priority preemptive scheduling is used for scheduling between the different queues.**

seeSSSS
Let USE° oskler "scheduling algorithm. P "M
queues are written below in the order of
assigned to them. V
priority and 5 is lowest priority.

1. System processes queue

2. Interactive processes queue

3. Interactive editing processes queue

4. Batch processes queue

5. Student processes queue

higher priority queue has complete priority.

Even if one process is completed and until higher priority queue becomes empty, processes in the queue cannot run.

For example, processes in the interactive editing queue could not execute unless and until the system processes in the queue finish the execution.

Processes in the interactive editing queue becomes empty.

If process in lower priority queue is executed the same time other process belonging to the same priority queue arrives then currently executing process

in lower priority queue should be preempted. If, a process entered the ready queue while a student process was executing, the student process would be preempted and batch process would be scheduled for execution.

In this algorithm, a fixed amount of CPU time is allocated to each queue and within this time different processes from this queue are scheduled.

For example, the interactive process queue can be given 70 percent of the CPU time for round robin scheduling among its processes, whereas the batch queue get 30 percent of the CPU time to give to its processes on an FCFS basis.

Highest priority

System processes

Interactive processes

Interactive editing processes

Batch process

Student processes

Lowest priority

→ 2.16.2(F) Multilevel Feedback-Queue Scheduling

Q- Explain multilevel feedback queue scheduling algorithm.

- In multilevel queue scheduling processes are not permitted to transfer from one queue to the other. Also queue assigned to the processes are permanent and cannot be changed, leading to the low scheduling cost, but it is not flexible.
- The multilevel feedback-queue scheduling algorithm processes can move from one queue to other. The CPU-bound processes use more CPU time and hence it will be transferred to a lower-priority queue.
- I/O-bound processes use less CPU time. Keeping CPU bound processes in low priority queue automatically leads to keeping the I/O bound and interactive processes in the higher-priority queues.
- There is a chance of starvation because of processes waiting for longer period of time in a lower-priority queue. It is avoided by aging; by transferring these waiting processes in a higher-priority queue.
- Consider the example of a multilevel feedback-queue scheduler with four queues, queue-0, queue-1, queue-2 and queue-3. Initially the scheduler starts executing all processes in queue-0.
- Once all the processes in queue-0 finish the execution, and queue-0 becomes empty, scheduler will consider the processes in queue-1 for execution. After queue-1 becomes empty then scheduler will consider the processes in queue-2 for execution.
- In the same way, processes in queue-3 will only be executed if queue-0, queue-1 and queue-2 are empty.
- A process arriving for queue-1 will preempt a process in queue-2. In the same way process that arrives for queue-0 will preempt a process in queue-1.
- Similarly if process arrives for queue-2 preempts process in queue-3. A process from the ready queue is placed in queue-0.

Multilevel feedback queue scheduling with 3 queues are shown in figure 3.9. If the process does not finish the execution within allocated time quantum of in queue-0, it is moved to the tail of queue-1.

The process at the head of queue-1 is given a quantum of 16 milliseconds if and only queue-0 is empty.

Process Management

- If it does not finish within the allocated time quantum, the process is placed into queue-2, and queue-2 are run on an FCFS basis but are run only if queue-0 and queue-1 are empty.
- If the process takes less than this time quantum, it gives highest priority to it taking 8ms gives highest CPU burst. It then moves to the tail of queue-1 and 8ms off its next burst.
- Processes in queue-2 are also served below 24 milliseconds. Lower priority processes with shorter burst times automatically go down to queue-2 and are served in FCFS order with any CPU cycles left over from queues 0 and 1.

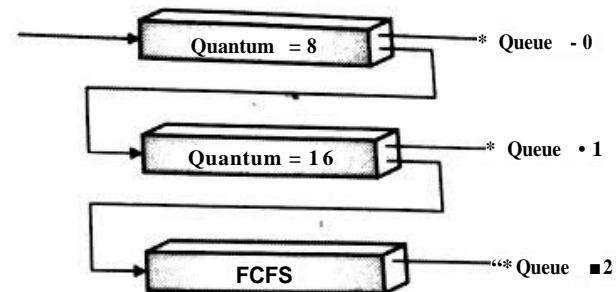


Fig. 2.16.2 : Multilevel feedback queues

2.17 Examples on Uniprocessor Scheduling Algorithms

Example 2.17.1 MU - Dec. 2014, 10 Marks

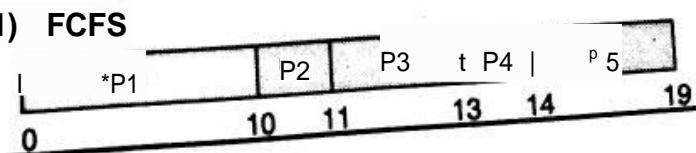
Consider the following set of processes, with the length of the CPU burst given in milliseconds:

Process	Burst Time	Priority
P1	10	3
P2	1	2
P3	2	4
P4	1	1
P5	5	2

The processes are assumed to have arrived in the order P1, P2, P3, P4, P5 at time 0. Draw a Gantt chart to illustrate the execution of these processes using the following scheduling algorithms: FCFS, WFQ, and RR (quantum = 1) and also calculate turnaround and average waiting time.

Solution:

(1) FCFS



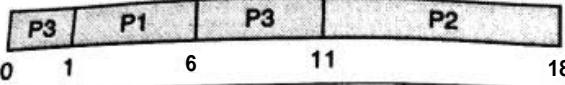
Process	Priority	TAT	Waiting time
P1	I	(6 - 1) = 5	(1 - 0) = 0
P2	3	(18 - 5) = 13	(11 - 5) = 6
P3	2	(11 - 0) = 11	(6 - 1) = 5
Average	j	9.6	3.66

$$\text{Average TAT} = (5 + 13 + 11)/3 = 9.6$$

$$\text{Average waiting time} = (0 + 6 + 5)/3 = 3.66$$

(ii) Preemptive priority

Following is the Gantt chart



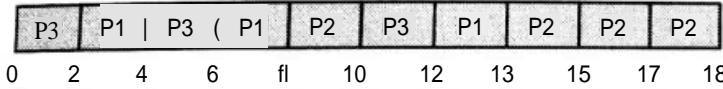
Process	Priority	TAT	Waiting time
P1	1	(6 - 1) = 5	(1 - 1) = 0
P2	3	(18 - 5) = 13	(11 - 5) = 6
P3	2	(11 - 0) = 11	(6 - 1) = 5
Average	j	9.66	3.66

$$\text{Average TAT} = (5 + 13 + 11)/3 = 9.6$$

$$\text{Average waiting time} = (0 + 6 + 5)/3 = 3.66$$

(iii) RR

Following is the Gantt chart



Process	Priority	TAT	Waiting time
P1	1	(13 - 1) = 12	(2 - 1) + (6 - 4) + (12 - 8) = 7
P2	3	(18 - 5) = 13	(8 - 5) + (13 - 10) = 6
P3	2	(12 - 0) = 12	(4 - 2) + (10 - 6) = 6
Average	j	12.33	6.33

$$\text{Average TAT} = (12 + 13 + 12)/3 = 12.33$$

$$\text{Average waiting time} = (7 + 6 + 6)/3 = 6.33$$

Example 2.17,3

Consider the four processes P1, P2, P3 and P4 with length of CPU burst time. Find out average waiting time and average turnaround time for the following algorithm.

(i) FCFS

00 RR (slice = 4 ms)

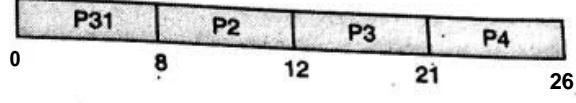
W SJF

Process	Arrival time	Burst time
P1	0	8
P2	1	4
P3	2	9
P4	3	5

Solution :

(0) FCFS

Following is the Gantt chart



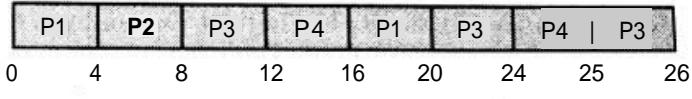
Process	TAT	Waiting time
P1	(8 - 0) = 8	0
P2	(12 - 1) = 11	(8 - 1) = 7
P3	(21 - 2) = 19	(12 - 2) = 10
P4	(26 - 3) = 23	(21 - 3) = 18
Average	15.25	8.75

$$\text{Average TAT} = (8 + 11 + 19 + 23)/4 = 15.75 \text{ ms}$$

$$\text{Average waiting time} = (0 + 7 + 10 + 18)/4 = 8.75 \text{ ms}$$

(ii) RR (slice = 4 ms)

Following is the Gantt chart



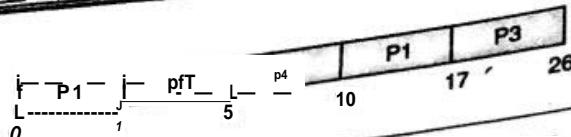
Process	TAT	Waiting time
P1	20	(16 - 4) = 12
P2	(8 - 1) = 7	(4 - 1) = 3
P3	(26 - 2) = 24	(8 - 2) + (20 - 12) + (25 - 24) = 15
P4	(25 - 3) = 22	(12 - 3) + (24 - 16) = 17
Average	18.25	11.75

$$\text{Average TAT} = (20 + 7 + 24 + 22)/4 = 18.25 \text{ ms}$$

$$\text{Average waiting time} = (12 + 3 + 15 + 17)/4 = 11.75 \text{ ms}$$

(iii) Preemptive SJF

Following is the Gantt chart



Process	tat	Waiting time
P1	(17 - 0) = 17	(10 - 1) * 9
P2	(5 - 1) = 4	(1 - 0) = 0
P3	(26 - 2) = 24	07 2 = 15
P4	(10 - 3) = 7	<5 - 3) = 2
Average	13	6.5

$$\text{Average TAT} = (17 + 4 + 24 + 7)/4 = 13 \text{ nis.}$$

Average waiting time

Example 2.17.4 for execution at suppose mat the following process amve time indicate

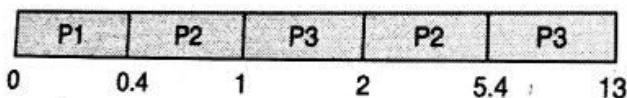
Process	Arrival time	Burst time
P1	0,0	8
P2	0.4	4
P3	1.0	1

Calculate average waiting time and average turnaround time using SJF and SRTF.

Solution:

(i) Preemptive SJF is SRTF (Shortest Remaining Time Next)

Following is the Gantt chart



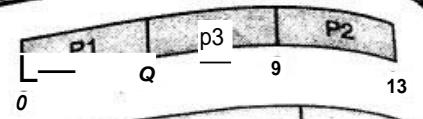
Process	TAT	Waiting time
P1	(13 - 0) = 13	(5.4 - 0.4) = 5
P2	(5.4 - 0.4) = 5	(2 - 1) = 1
P3	(2 - 1) = 1	(1 - 1) = 0
Average	6.33	2

$$\text{Average TAT} = (13 + 5 + 1)/3 = 6.33$$

$$\text{Average waiting time} = (5 + 1 + 0)/3 = 2$$

W Non-preemptive SJF

Following is the Gantt chart



Process	TAT	Waiting time
P1	(8 - 0) = 8	0
P2	(12 - 5) = 7	(9 - 0.4) = 8.6
P3	(8 - 1) = 7	(8 - 1) = 7
Average	9.53	5.2

$$\text{Average TAT} = (8 + 12 + 8)/3 = 9.53$$

$$\text{Average waiting time} = (0 + 8.6 + 7)/3 = 5.2$$

Example 2.17.5

consider the following set of processes having the burst time

Process	CPU Burst time	Arrival time
P1	10	0
P2	5	1
P3	2	2

Calculate average waiting time using the following scheduling algorithms.

(D) FCFS

(2) SJF

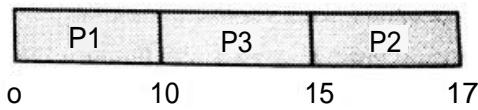
(3) Priority scheduling having priority range from 1 to 3 respectively for process P1 = 3, P2 = 2, P3 = 3 as

(4) RR (slice = 2)

Solution :

(1) FCFS

Following is the Gantt chart

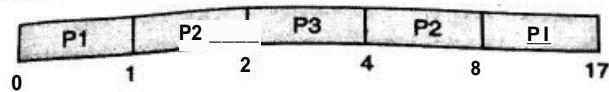


Process	Waiting time
P1	0
P2	(10 - 1) = 9
P3	(15 - 2) = 13
Average	7.33

$$\text{Average waiting time} = (0 + 9 + 13)/3 = 7.33$$

(2) Preemptive SJF

Following is the Gantt chart

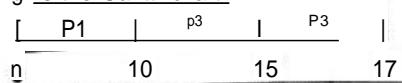


process	Waiting time
P1	$(8 - 1) = 7$
P2	$(4 - 2 - 1) = 1$
P3	$(2 - 3) = 0$
Average	2.67

$$\text{Average waiting time} = (7 + 1 + 0)/3 = 2.67 \text{ ms.}$$

Priority scheduling having priority range from 1 to 3. respectively for process P1 = 3, P2 = 2, P3 = 3 as given.

Following is the Gantt chart

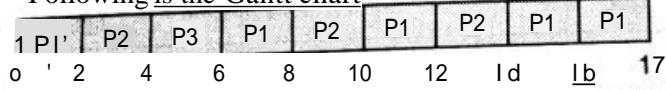


Process	Waiting time
P1	0
P2	$(10 - 1) = 9$
P3	$(15 - 2) = 13$
Average	7.33

$$\begin{aligned} \text{Average waiting time} &= (0 + 9 + 13)/3 \\ &= 7.33 \text{ ms. (Same as FCFS)} \end{aligned}$$

(4) RR (slice = 2)

Following is the Gantt chart



Process	Waiting time
P1	$(6 - 2) + (10 - 8) + (13 - 12) = 7$
P2	$(2 - 0) + (8 - 4) + (12 - 10) = 7$
P3	$(4 - 2) = 2$
Average	5.33

$$\text{Average waiting time} = (7 + 7 + 2)/3 = 5.33 \text{ ms.}$$

Example 2.17.6

Consider the following set of processes.

Process	CPU Burst time	Arrival time
P1	3	0
P2	5	1
P3	2	2
P4	5	3
P5	5	4

Calculate average waiting and turnaround time for each algorithm.

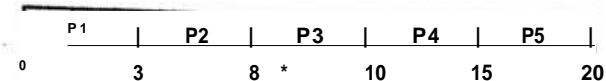
(1) FCFS

(2) SJF

(3) RR (slice = 2)

Solution :

(1) FCFS : Following is the Gantt chart

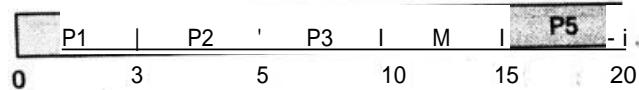


Process	TAT	Waiting time
P1	$(3 - 3) = 0$	0
P2	$(8 - 1) = 7$	$(3 - 1) = 2$
P3	$(10 - 2) = 8$	$(8 - 2) = 6$
P4	$(15 - 3) = 12$	$(10 - 3) = 7$
P5	$(20 - 4) = 16$	$(15 - 4) = 11$
Average	9.2	5.2

$$\begin{aligned} \text{Average turnaround time} &= (0 + 7 + 8 + 12 + 16)/5 \\ &= 9.2 \text{ ms.} \end{aligned}$$

$$\begin{aligned} \text{Average waiting time} &= (0 + 2 + 6 + 6 + 11)/5 \\ &= 5.2 \text{ ms.} \end{aligned}$$

(2) SJF : Following is the Gantt chart



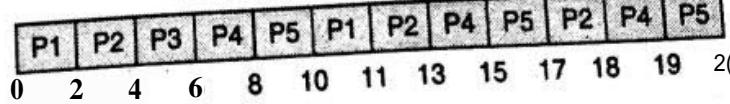
Process	TAT	Waiting time
P1	$(3 - 3) = 0$	0
P2	$(10 - 1) = 9$	$(5 - 1) = 4$
P3	$(5 - 2) = 3$	$(3 - 2) = 1$
P4	$(15 - 3) = 12$	$(10 - 3) = 7$
P5	$(20 - 4) = 16$	$(15 - 4) = 11$
Average	8.6	4.6

$$\begin{aligned} \text{Average turnaround time} &= (0 + 9 + 3 + 12 + 16)/5 \\ &= 8.6 \text{ ms.} \end{aligned}$$

$$\begin{aligned} \text{Average waiting time} &= (0 + 4 + 1 + 7 + 11)/5 \\ &= 4.6 \text{ ms.} \end{aligned}$$

(3) RR (slice = 2)

Following is the Gantt chart



Process	TAT	Waiting time
P1	$(11 - 0) = 11$	$(0 - 0) = 0$
P2	$(18 - 1) = 17$	$(2 - 1) + (11 - 4) + (17 - 13) = 2$
P3	$(6 - 2) = 4$	$(4 - 2) = 2$
P4	$(19 - 3) = 16$	$(6 - 1) + (13 - 8) + (18 - 15) = 11$
P5	J20-4KJ1	$(8 - 4) + (15 - 10) + (19 - 17) = 11$
Average	12.8	8.8

$$\begin{aligned}
 \text{Average turnaround time} &= \frac{(11 + 17 + 4 + 16 + 11)}{5} \\
 &= 12.8 \text{ ms} \\
 &= \frac{(8 + 12 + 2 + 11 + 11)}{5} \\
 &= 8.8 \text{ ms.}
 \end{aligned}$$

Example 2.17.7

Assume that you have 5 jobs to execute with one processor

Job	CPU Burst time	Arrival time
0	75	0
1	50	10
2	25	10
3	20	80
4	45	85

Suppose system uses round robin with quantum of 15.

- Draw Gantt chart
- Find average wait and turnaround time

Soln. : The Gantt chart is,

Process	Priority	Arrival time	Completion time	TAT	Waiting time
P1	3	0	8	$(8 - 0) = 8$	0
P2	1	10	15	$(15 - 10) = 5$	5
P3	2	10	17.5	$(17.5 - 10) = 7.5$	7.5
P4	3	80	100	$(100 - 80) = 20$	20
P5	4	85	130	$(130 - 85) = 45$	45
Average				$(8 + 5 + 7.5 + 20 + 45) / 5 = 130 / 5 = 26$	26

Process	TAT	Waiting time
P0	$(215 - 0) = 215$	$(45 - 15) + (115 - 60) + (165 - 30) + (200 - 180) = 140$
P1	$(185 - 10) - 175$	$(15 - 10) + (60 - 30) + (130 - 75) + (180 - 145) = 125$
P2	$(85 - 10) = 75$	$(30 - 10) + (75 - 45) = 50$
P3	$(150 - 50) = 100$	$(85 - 80) + (145 - 100) = 50$
P4	$(200 - 85) = 115$	$(100 - 85) + (150 - 115) + (185 - 165) = 70$
Average	130	87

$$\begin{aligned}
 \text{Average turnaround time} &= \frac{(215 + 175 + 75 + 70 + 115)}{5} \\
 &= 130 \text{ ms.}
 \end{aligned}$$

$$\begin{aligned}
 \text{Average waiting time} &= \frac{(14Q + 12)}{5} = \frac{(14 \times 5 + 12)}{5} = 12.8 \text{ ms.} \\
 &= 87 \text{ ms.}
 \end{aligned}$$

2.17.8
Use following Schedule for the following process

- Nonpreemptive SJF
- Preemptive SJF

Preemptive priority.

	Arrival time	Burst time	Priority
P1	0	8	3
P2	1	1	1
P3	2	3	2
P4	3	2	3
P5	4	6	4

Solution

- Round robin

Process	Priority	Arrival time	Completion time	TAT	Waiting time
P1	3	0	8	$(8 - 0) = 8$	0
P2	1	10	15	$(15 - 10) = 5$	$(15 - 8) = 7$
P3	2	10	17.5	$(17.5 - 10) = 7.5$	$(17.5 - 15) = 2.5$
P4	3	80	100	$(100 - 80) = 20$	$(100 - 17.5) = 82.5$
P5	4	85	130	$(130 - 85) = 45$	$(130 - 100) = 30$
Average				$(8 + 5 + 7.5 + 20 + 45) / 5 = 130 / 5 = 26$	$(8 + 7 + 2.5 + 82.5 + 30) / 5 = 130 / 5 = 26$

(ii) SJF(Preemptive)

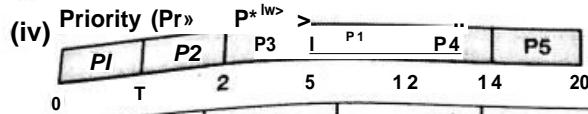
Process	Priority	Arrival time	Completion time	TAT	Waiting time
P1	3	0	8	$(8 - 0) = 8$	0
P2	1	10	11	$(11 - 10) = 1$	$(11 - 8) = 3$
P3	2	10	13	$(13 - 10) = 3$	$(13 - 11) = 2$
P4	3	80	82	$(82 - 80) = 2$	$(82 - 13) = 69$
P5	4	85	91	$(91 - 85) = 6$	$(91 - 82) = 9$
Average				$(8 + 1 + 3 + 2 + 6) / 5 = 130 / 5 = 26$	$(8 + 3 + 2 + 69 + 9) / 5 = 130 / 5 = 26$

Process	Priority	Arrival time	Completion time	TAT	Waiting time
P1	3	0	8	$(8 - 0) = 8$	0
P2	1	10	11	$(11 - 10) = 1$	$(11 - 8) = 3$
P3	2	10	13	$(13 - 10) = 3$	$(13 - 11) = 2$
P4	3	80	82	$(82 - 80) = 2$	$(82 - 13) = 69$
P5	4	85	91	$(91 - 85) = 6$	$(91 - 82) = 9$
Average				$(8 + 1 + 3 + 2 + 6) / 5 = 130 / 5 = 26$	$(8 + 3 + 2 + 69 + 9) / 5 = 130 / 5 = 26$

(iii) SJF(Nonpreemptive)

Process	Priority	Arrival time	Completion time	TAT	Waiting time
P1	3	0	8	$(8 - 0) = 8$	0
P2	1	10	11	$(11 - 10) = 1$	$(11 - 8) = 3$
P3	2	10	13	$(13 - 10) = 3$	$(13 - 11) = 2$
P4	3	80	82	$(82 - 80) = 2$	$(82 - 13) = 69$
P5	4	85	91	$(91 - 85) = 6$	$(91 - 82) = 9$
Average				$(8 + 1 + 3 + 2 + 6) / 5 = 130 / 5 = 26$	$(8 + 3 + 2 + 69 + 9) / 5 = 130 / 5 = 26$

Proc	Priority	Finish time	TAT	Waiting time
P1	3	8	(8-0) - B	0
P2	1	9	(9-1) - a	(8-1)*7
P3	2	12	(14-2) = 12	(11-2) = 9
P4	3	14	(11-3) = 8	(9-3) = 6
PS	4	20	(20-4) = 16	(14-4) = 10
Average			(52/5) = 10.4	(32/5) = 6.4



Proc	Priority	Finish time	TAT	Waiting time
P1	3	20	(12-0) = 12	(5-1) = 4
P2	1	2	(2-1) = 1	(1-1) = 0
P3	2	5	(5-2) = 3	(2-1) = 1
P4	3	7	(14-3) = 11	(12-3) = 9
P5	4	13	(20-4) = 16	(14-4) = 10
Average			(43/5) = 0.6	(23/5) = 4.6

Example 2J 7.9

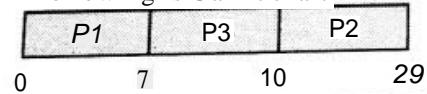
Find AWT, ATAT, ART and AWTAT for the following set of processes with CPU burst time in ms. Assume that all processes arrive at time 0.

(P1-19), (P2-7), (P3-3)

- FCFS with order P2, P3, P1
- Round Robin (Quantum = 2ms)

Solution :

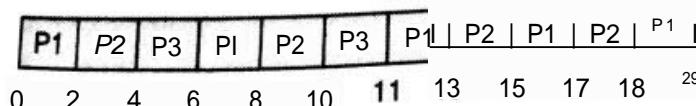
(i) FCFS : Following is Gantt chart



Process	TAT	Waiting time
P1	(29-0) = 29	(10-0) = 10
P2	(7-0) = 7	0
P3	(10-0) = 10	(7-0) = 7
Average	46/3 = 15.33	17/3 = 5.66

(ii) Round Robin (RR) with quantum = 2

Following is Gantt chart



Process Managsm

Proc+ca	TAT	Waiting time
P1	(29-0) = 29	(B=2) + (11-B) + (16-13) + (16-17) = 10
P2	(11-0) = 11	(2-0) + (6-4) + (13-10) + (17-15) = 11
P3	(13-0) = 13	(4-0) + (10-6) > B
Average	W3 = 18.13	129/3 = 9.66

Example 2.17.10 MU - Dre. 2016. 10 Marks

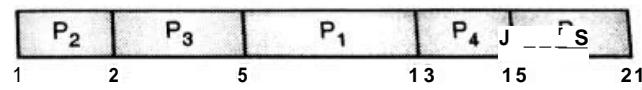
Assume the following processes arrive for execution at the time indicated and the length of each burst time given in msec.

Job	Burst time	Priority	Arrival time
P1	8	3	3
P2	1	1	A
P3	3	2	?
P4	2	3	3
P5	6	4	*

For the above process parameters, find average waiting times and average turnaround times for the following scheduling algorithms. First Come First Serve, Shortest Job First, non preemptive priority and Round Robin (assume quantum = 2 Units)

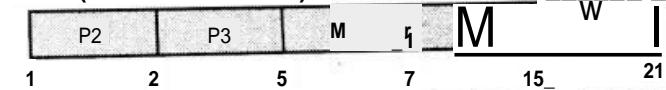
Solution :

FCFS (First Solution)



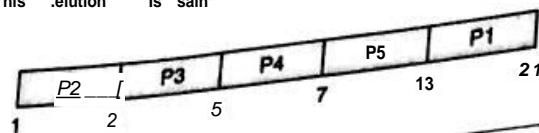
Job	Priority	Arrival Time	Finish time	TAT	Waiting time
P1	3	3	13	(13-3) = 10	(53) = 2
P2	1	1	2	(2-1) = 1	(1-1) = 0
P3	2	2	5	(5-2) = 3	(2-2) = 0
P4	3	3	15	(15-3) = 12	(13-3) = 0
P5	4	4	21	(21-4) = 17	(15-4) = 11
Average				(43/5) = 8.6	(23/5) = 4.6

FCFS (Second Solution)



Job	Priority	Arrival Time	Finish time	TAT	Waiting time
P1	3	3	15	(15-3) = 12	(P-3) = 4
P2	1	1	2	(2-1) = 1	(1-1) = 0
P3	2	2	5	(5-2) = 3	(2-2) = 0
P4	3	3	7	(7-3) = 4	(5-3) = 2
P5	4	4	21	(21-4) = 17	(15-4) = 11
Average				(37/5) = 7.4	(17/5) = 3.4

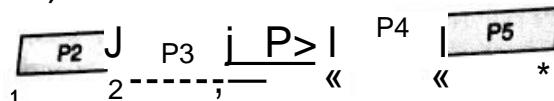
No. (This duration is same for SJF preemptive)



Job	Priority	Arrival Time	Finish Time	TAT	Waiting time
P1	3	3	21	(2f4) = 18	(13-3) = 10
P2	1	1	2	(2-1) = 1	(1-1) = 0
P3	2	2	5	(S2)-3	(2-2) = 0
P4	3	3	7	(7-3)-4	(5-3) = 2
P5	4	4	13	(13-4)-9	(74) = 3

Nonpreemptive Priority: (Same as SJF)

FCFS



Job	Priority	Arrival Time	Finish time	TAT	Waiting time
P1	3	3	13	(13-3) = 10	(5-3) = 2
P2	1	1	2	(2-1) = 1	(1-1) = 0
P3	2	2	5	(5-2) = 3	(2-2) = 0
P4	3	3	15	(15-3) = 12	(13-3) = 10
P5	4	4	21	(21-4) = 17	(15-4) = 11
Average				(43/5) = 8.6	(23/5) = 4.6

Round Robin - quantum - 2 units

Job	Priority	Arrival Time	Finish time	TAT	Waiting time
P1	3	3	21	(21-3) = 18	(4-3) + (11-6) + (15-13) + (19-15) = 12
P2	1	1	2	(2-1) = 1	(M) 0
P3	2	2	11	(11*2) = 9	(2-2) + (104) = 6
P4	3	3	2	(3-3) = 0	(6-3) = 3
P5	4	4	19	(19-4) = 15	(8-4) + (13-10) + (17-15) = 9
Average				(48/5) = 9.6	(005) = 6

Example 2.17.11

Consider following set of processes with given priorities

Process	Burst Time*	Arrival Time	Priority
P1	8	0	3
P2	1	1	1
P3	3	2	2
P4	2	3	3
P5	6	4	4

Draw Gantt chart for FCFS, SJF, Preemptive and RR (Round Robin) algorithms

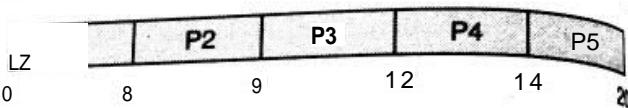
What is the turnaround time of each process for above algorithms?

What is waiting time of each process for each of the above algorithms?

Which algorithm results in minimum average waiting time?

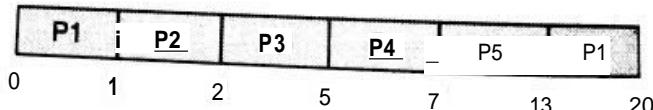
Solution :

FIFO



Process	Priority	Finish time	TAT	Waiting time
P1	3	8	(8-0) = 8	0
P2	1	9	(9-1) = 8	(8-1) = 7
P3	2	12	(12-2) = 10	(9-2) = 7
P4	2	14	(14-3) = 11	(12-3) = 9
P5	4	20	(20-4) = 16	(14-4) = 10
Average			(53/5) = 10.6	(33/5) = 6.6

SJF (Preemptive)

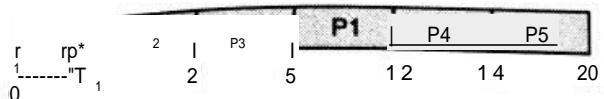


Process	Priority	Finish time	TAT	Waiting time
P1	3	20	(20-0) = 20	(13-1) = 12
P2	1	2	(2-1) = 1	(1-1) = 0
P3	2	5	(5-2) = 3	(2-2) = 0
P4	3	7	(7-3) = 4	(5-3) = 2
P5	4	13	(13-4) = 9	(7-4) = 3
Average				

JjHN nP (tive)

	P1	P2	P4	P3	P5	
	0	8	9	11	14	20
Process	ftforW	Finish time	TAT	Waiting time		
P1	3	8	(8-0) = 8	0		
P2	1	9	(9-1) = 8	(8-1) = 7		
P3	2	12	(14-2) = 12	(11-2) = 9		
P4	3	14	(11-3) = 8	(8-3) = 5		
P5	4	20	(20-4) = 16	(14-4) = 10		
Avg			(52/5) = 10.4	(32/5) = 6.4		

priority(Preempt)



Procs*	Priority	Finish time	TAT	Waiting time
P1	3	20	(12-0) = 12	(5-1) = 4
P2	1	2	(2-1) = 1	(1-1) = 0
P3	2	5	(5-2) = 3	(2-2) = 0
P4	3	7	(14-3) = 11	(12-3) = 9
P5	4	13	(20-4) = 16	(14-4) = 10
Average			(43/5) = 8.6	(23/5) = 4.6

RR(Quantum 2)



Process	Priority	Finish time	TAT	Waiting time
P1	3	20	(20-0) = 20	0+(9'2)+(14-1)H18-16) = 12
P2	1	3	(3-1) = 2	(2-0) = 1
P3	2	12	(12-2) = 10	(3-2)M11-5) = 7
P4	3	7	(7-3) = 4	(5-3) = 2
P5	4	18	(18-4) = 14	(7-4)+(12-9)+(16-14) = 8
Average			(50/5) = 10	(30/5) = 6

lives minimum average waiting lieie.

Preemptive SJF g

Example 2.17.12 MU - Nov. 2015. 10 Marks

Consider following set of processes with their CPU burst time.

Process	Burst time	Arrival time
P1	10	1
P2	4	2
P3	5	3
P4	3	4

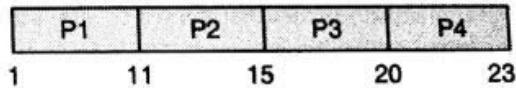
(i) Draw Gantt chart FCFS, SJF preemptive and round robin (quantum w 3), Calculate average waiting time and average turnaround time.

(ii) Explain which scheduling policy adopted by Linux?

Solution :

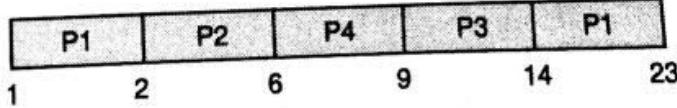
Process	TAT	Waiting time
P1	(23-1) = 22	(14-2) = 12
P2	(6-2) = 4	0
P3	(14-3) = 11	(9-3) = 6
P4	(4-4) = 0	(6-4) = 2
Average	42/4 = 10.5	20/4 = 5

(i) FCFS Scheduling

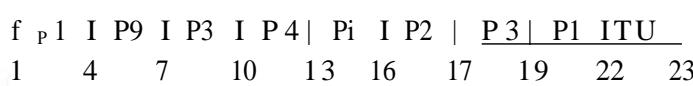


Process	TAT	Waiting time
P1	(11-1) = 10	0
P2	(15-2) = 13	(11-2) = 9
P3	(20-3) = 17	(15-3) = 12
P4	(23-4) = 19	(20-4) = 16
Average	59/4 = 14.75	37/4 = 9.25

(H) SJF Preemptive Scheduling



(iii) RR Scheduling with quantum = 3





Process	TAT	Waiting time
P1	$(23-1)=22$	$(13-4)+(19-16)=12$
P2	$(17-2)=15$	$(4-2)+(16-7)=11$
P3	$(19-3)=16$	$(7-3)+(17-10)=11$
P4	$(13-4)=9$	$(10-4)=6$
Average	$62/4=15.5$	$40/4=10$

Syllabus Topic : Thread Scheduling**2.78 Thread Scheduling****2.18.7 Contention Scope**

Q. Write note on thread scheduling

- Scheduling of the user level and kernel level threads is different.
- If system implements many to one or many to many model of multithreading then thread library schedules user-level threads to run on an available LWP.
- This scheme of scheduling is called as process-contention scope (PCS) as threads belonging to the same process competes for processor.
- The operating system has to actually schedule the kernel thread on physical processor to run a user thread scheduled by thread library on available LWPs.
- Kernel uses system-contention scope (SCS) to schedule the threads on processor when there are many threads competing for the physical processor.
- Windows XP, Solaris 9; and Linux use one-to-one model and schedule threads using only system-contention scope.

Usually, PCS is carried out in keeping with the scheduler chooses the runnable thread with the highest priority to run.

Programmers set the priorities

which are not altered by the system. Some thread libraries may permit the

alter the priority of a thread.

The system will preempt the thread present for a higher-priority thread. On the other hand, there is no assurance of time sharing among threads of equal priority.

2.18.2 pthread Scheduling

- Following are the pthread value, PTHREAD_SCOPE_PROCESS is used to threads using PCS scheduling. PTHREAD_SCOPE_SYSTEM is used to threads using SCS scheduling.

The system that implements many PTHREAD_SCOPE_PROCESS to schedule available light weight processes. Thread library the number of LWPs. Scheduler activates this purpose. The PTHREAD_SCOPE and bind an LWP for each user-level many systems, successfully mapping threads the one-to-one policy.

In order to get and set contention Pthread IPC offers following two functions

- pthread_attr_setscope(pthread_attr_t *attr, int scope)
- pthread_attr_getscope(pthread_attr_t *attr, int *scope)

Syllabus Topic : Multiple Processor Scheduling**2.19 Multiple-Processor Scheduling**

Q. Explain multiprocessor scheduling

Here we assume the homogeneous multiprocessors are identical.

2.19.1 Approaches to Multiple-Processor Scheduling

- In asymmetric multiprocessing, all scheduling decisions, I/O processing, and other system tasks are carried out by master processor.
- All other processors execute in user mode. System data structures are accessed by single processor leading to minimizing the data sharing.
- In symmetric multiprocessing, scheduling is carried out by each processor.
- **Preemptive Scheduling:** The processor can maintain a queue for ready processes. Processor may have its own queue for ready processes.
- All modern operating system supports symmetric multiprocessing.

2.19.2 Processor Affinity

- Some systems also use affinity scheduling in which same thread again and again gets scheduled on same CPU. This makes better use of previously cached blocks and increases the cache hits.
- Also, the TLB may also have the right pages, reducing TLB faults. For creation of affinity a two-level scheduling algorithm is used.
- After a thread is created, it is assigned to a lightly loaded CPU which is a top level of the algorithm. It ensures that, each CPU acquires its own collection of threads.
- The bottom level of the algorithm performs real scheduling. It is carried out by each CPU separately, on the basis of priorities or some other means.
- By trying to keep a thread on the same CPU for its entire lifetime, cache affinity is maximized.
- In case if a CPU has no threads to run, it takes one from another CPU instead of remaining idle.
- Approximately equal load distribution over the available CPUs, cache affinity and minimizing contention for the ready lists as CPU frequently uses same ready list are three benefits offered by two-level scheduling algorithm.
- Soft affinity means same process will be scheduled on same processor and keep running on it. This can be a policy but cannot be guaranteed by OS.
- Process migration on other CPU may happen. In hard affinity process can not be migrated to other processor.

2.19.3 Load Balancing

- The main objective of load balancing in symmetric multiprocessing systems (SMP) is to keep the load balanced among all the processors.
- Load is evenly distributed among all processors. Otherwise some processors may remain idle and others will be heavily loaded.
- In push migration approach of load balancing, a specific process periodically checks load of processors. If it finds the processor with light load then it migrates process from heavily loaded processor to lightly loaded or idle processor.
- In pull migration approach, idle processor pulls waiting task from heavily loaded processor.

2.19.4 Multicore Processors

- Now days, it is possible to place multiple processor cores on the same physical chip.
- This is called multicore processor. Each of the cores is treated as separate physical processor by OS as there is a register in core to set its architectural state.
- SMP systems using multicore processors consume less power and are efficient.
- Scheduling is complex for multicore system. The processor accessing memory has to wait for data as it is not available. So processor spends 50% of its time in waiting.
- As a solution to this, hardware designer provided two or many hardware threads for each core.
- If one thread has to wait for memory then core switches to other thread. OS treats each hardware thread as a logical processor for S/W thread.
- The coarse grained and fine grained multithreading can be used to multithread a processor.
- Multithreaded multicore processor requires two levels of scheduling. OS schedules software thread on hardware thread using any scheduling algorithm.
- Second level of scheduling involves how each core makes a decision about which hardware thread to execute.

2.19.5 Virtualization and Scheduling

- With virtualization, single CPU system works like multiprocessor system.
- The virtualization software offers one or more virtual processors to each of the VMs running on the system and then schedules the use of the real CPUs among the VMs. VMs are created and managed by host operating system. Each VM has guest OS and applications running within it.
- The guest OS scheduling algorithm that expects a convinced amount of progress in a given amount of time will be negatively impacted by virtualization.

2.19.6 Other Multiprocessor Scheduling Approaches

Following are the four general approaches used for scheduling of threads on multiprocessor system.

1. Is the

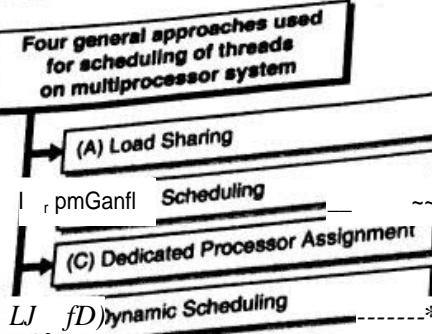


Fig. C.2.12: Approaches for scheduling of threads on multiprocessor system

→ (A) Load Sharing

All the ready threads are present in single global queue and each processor selects thread from this queue when it is idle.

→ (B) Gang Scheduling

A set of related threads is scheduled to execute on a set of processors all at once, on a one-to-one basis.

→ (C) Dedicated Processor Assignment

In this approach, during execution of program the number of processors assigned is equal to number of threads available in program. All the processors return to pool of processors after program complete the execution.

→ (D) Dynamic Scheduling

The number of threads in process may changes during the execution of that process.

2.19.6(A) Load Sharing

Advantages

- In this approach, load is equally distributed among all the processors so that no one remains idle.
- There is no need of centralized scheduler. The scheduling routine of the operating system executes on available processor to choose the next thread.

The maintained global queue can be accessed by using schemes used in uniprocessor scheduling such as

Round robin schemes or schemes based on execution

Load sharing approaches

Load sharing approaches

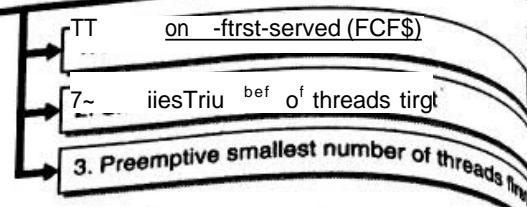


Fig. C.2.13. Load sharing approaches

→ 1. Non-preemptive first-come-first-served (FCFS)

The threads of the newly arrived job are placed consecutively at the end of shared queue. The processor selects next ready thread and executes until it finishes execution or gets blocked.

→ 2. Smallest number of threads first

If the jobs are having smallest number of threads then highest priority is given to it. The ready queue is like priority queue and jobs are ordered according to jobs arrived first. To FCFS, thread executes until it finishes execution or gets blocked.

→ 3. Preemptive smallest number of threads first

Job having smallest number of threads gets highest priority. If an arriving job is having same number of threads with compare to executing job it preempts the threads of scheduled job.

Disadvantages

- A single shared ready queue may occupy the memory portion and if many processors simultaneously access it then it may become bottleneck.
- Preempted threads may not get same processor to resume execution. The caching becomes less effective as processors have cache memory.
- It is unlikely to get access to all the threads in program at the same time. If more coordination is required between threads then process switch may lead to poor performance.

2.19.6(B) Gang Scheduling

Advantages

- As synchronization, blocking can be minimized. Related processes can be achieved by running parallel threads which increases performance.

Scheduling overhead can be minimized as a single decision affects a number of processors and processes at one time.

Group scheduling is better for the application when an application whose any part is not running and other is ready for execution. It is implemented in many multiprocessor operating

systems. In this scheduling process switches are multiplexed and rejected threads are ran in parallel. Hence, performance is good.

As threads that require coordination run in parallel, they can access file without additional overhead and resource allocation can be done with less overhead.

In this scheduling processor allocation is required. If N processors and M applications with N or fewer threads are present then each application could be given $1/M$ of the available time on the N processors, using time slicing.

→ 2.19.6(C) Dedicated Processor Assignment

- In this approach, during execution of program the number of processors assigned is equal to number of threads available in program.
- All the processors returns to pool of processors after program complete the execution. In this approach, if thread of an application is blocked waiting for I/O or for synchronization with other thread then processor remains idle. It can be addressed as following :
 - o If tens or hundreds of processors present in system then processor utilization does not matter for performance,
 - o As process switching is avoided, it results in better improvement in performance,

→ 2.19.6(D) Dynamic Scheduling

- In some applications, the number of threads changes during execution.
- The operating system then can adjust the load to improve performance. Both OS and application can carry out scheduling task.
- The operating system partitions the processors to allocate to the jobs,
- The runnable task of each job is mapped to threads and is given to the processors for execution at present in its partition.

A proper decision, about which subset to run, in addition to which thread to suspend when a process is

preempted, is left to the individual applications (maybe through a set of nm-time library routines).

- This approach may be appropriate for some applications and not for all. Applications can be implemented to take advantage of this feature of operating system.

Operating system is only responsible to processor allocation and performs following tasks when job demands one or more processors.

- Allocate the idle processor to satisfy the request.
- If job is new arrival and need processor then take the processor from currently allocated job to which more than one processor is allocated and allocate to newly arrived job.
- In case if job's request cannot be satisfied then it remains outstanding until processor becomes available or the job withdraws the request.

After release of one or more processors

- Search the queue of requests that are not satisfied for processors. Allocate a single processor to each job in the list that are waiting new arrivals and currently has no processors.
- Then search the list again, allocating the remaining processors on an FCFS basis.

2.20 Exam Pack (University and Review Questions)

☞ Syllabus Topic : Process Concept

- Q. What do you mean by process ?
(Refer section 2.2) (3 Marks) (Dec. 2014)
- Q. Explain the concept of process.
(Refer section 2.2)
- Q. What is context switch?
(Refer section 2.3)

☞ Syllabus Topic : Operation on Processes

- Q. What are the operations performed on process? (Refer section 2.4)
- Q. Explain process creation.
(Refer section 2.4.1)
- Q. Explain process termination
(Refer section 2.4.2)
- Q. Explain role of process control block.
(Refer section 2.5) (5 Marks) (Dec. 2014)

Q. Explain process control block.
(Refer section 2.5)

Q. Draw and explain process state transition diagram. (Refer section 2.6)

☛ **Syllabus Topic : Process Scheduling**

Q. Explain long-term scheduler.
(Refer section 2.8.1(A))

Q. Explain short-term scheduler.
(Refer section 2.8.1(B))

Q. Explain medium-term scheduler.
(Refer section 2.8.1(C))

Q. Compare the functions of different types of schedulers. (Refer section 2.8.1(D))

☛ **Syllabus Topic : Multithreading**

Q. Explain user level and kernel level threads with advantages and disadvantages.
(Refer section 2.11)

☛ **Syllabus Topic : Process - Multithreading Models**

Q. Explain various multithreading models.
(Refer section 2.12)

☛ **Syllabus Topic : Thread Libraries**

Q. Write note on thread libraries.
(Refer section 2.13)

☛ **Syllabus Topic : Threading Issues**

Q. What are different threading issues?
(Refer section 2.14)

☛ **Syllabus Topic : Process Scheduling : Basic Concepts**

Q. What is difference between preemptive and non-preemptive scheduling? (Refer section 2.15.2)

☛ **Syllabus Topic : Scheduling Criteria**

Q. What are the criteria for evaluating scheduling algorithm performance? (Refer section 2.16.1)

☛ **Syllabus Topic : Scheduling Algorithms**

Q. Explain FIFO scheduling algorithm.
(Refer section 2.16.2(A))

Q. Explain SJF and SRTN scheduling algorithm.
(Refer section 2.16.2(B))

Q. Describe priority scheduling algorithm.
(Refer section 2.16.2(C))

Q. With the help of example, explain Round Robin scheduling algorithm. (Refer section 2.16.2(D))

Q. Explain multilevel queue scheduling.
(Refer section 2.16.2(E))

Q. Explain multilevel feedback queue algorithm. (Refer section 2.16.2(F))

Example 2.17.1 (10 Marks)

Example 2.17.10 (10 Marks)

Example 2.17.12 (10 Marks)

☛ **Syllabus Topic : Thread Scheduling**

Q. Write note on thread scheduling.
(Refer section 2.18.1)

☛ **Syllabus Topic : Multiple Processor Scheduling**

Q. Explain multiprocessor scheduling in detail.
(Refer section 2.13)

CHAPTER 131

Process Coordination

[Module in]

Syllabus Topics

Synchronization : The critical Section P hi

SSZSS

Peterson's Solution, monitors, Atomic transaction, Methods for Recovery from Deadlock.

S - Z

Syllabus Topic : Synchronization

3.1 Background

→ (May 16)

Q. Explain the process synchronization in brief.

MU - May 2016, 5 Marks

- In a single processor multiprogramming system, processes are not executed concurrently. In order to get the appearance of concurrent execution, a fixed time slot is allocated to each process. After utilization of this time slot, CPU gets allocated to other process. Such switching of CPU back and forth between processes is called as context switch.
- At a time single process gets executed so parallel processing cannot be accomplished. Also there is a definite amount of overhead drawn in switching back and forth between processes. Apart from above limitations, interleaved execution offers major benefits in processing efficiency and in program structuring.
- In a multiple processor system, interleaving and overlapping the execution of multiple processes is achievable. Both interleaving and overlapping correspond to basically diverse modes of execution and present diverse problems. In reality, both interleaving and overlapping can be treated as illustration of concurrent processing and both the techniques address the similar problems.
- The comparative speed of execution of processes depends on activities and behavior of other processes.

how the operating system handles interrupts, and the scheduling policies of the operating system.

The difficulties that arise

1. Global resources sharing. For example, suppose two processes use the same global variable simultaneously and both carry out read and write operations on that variable, then various reads and writes execution ordering is serious.
2. Optimal management of the resources is not easy for the operating system.
3. Programming error tracing is not easy as results are typically non-deterministic and not reproducible.

Only a single user is supported by single processor multiprogramming system at a time. While working on one application, user can switch to other application in this system. The keyboard for input purpose and screen for output purpose used by each application is same. The reason is, each application needs to use the procedure echo. As echo procedure is shared by each application, memory can be saved by keeping single copy of procedure in memory portion which is global for all sharing applications.

Example

```
void echof()
{
    chinput = getchar();
    choutput = chinput;
    putchar(choutput);
}
```

In order to have efficient and close interaction among processes, sharing of memory among processes is useful. Consider the following sequence:

- P₁ calls the echo procedure. The initially process gets char returns its value and stores it in input variable chinput. Process P₁ gets interrupted. At this instant, the most recently entered character, m, is stored in chinput. At this instant, the most recently entered character, m, is stored in chinput. The process P₁ resumes its execution and the value of chinput is now again p. Now process P₂ calls the echo procedure. Process P₂ inputs and displays the single character n on monitor. After P₁ gets interrupted, there is turn of process P₂ and it call up the echo procedure. At this situation, process P₁ is in suspended state and it is still inside the echo procedure. The process P₂ is not allowed to enter inside the echo() procedure. Therefore, P₂ is suspended until P₁ comes out of the echo procedure. In this way, the first character m is lost and the second character n is displayed two times. All this happens due to shared global variable, chinput. If after updating the shared global variable, one process is interrupted, another process may modify the variable before the interrupted process can use its value. In above example both P₁ and P₂ are allowed to use shared global variable chinput. However, if only one process at a time is allowed to be in procedure, above discussed sequence would result in the following:
- Initially process P₁ calls the echo procedure. The getchar() returns its value and stores it in input variable chinput. Process P₁ gets interrupted straight away after returned value gets stored in chinput. At this instant, the most recently entered character, m, is stored in variable chinput.
 - After P₁ gets interrupted, there is turn of process P₂ and it call up the echo procedure. At this situation, process P₁ is in suspended state and it is still inside the echo procedure. The process P₂ is not allowed to enter inside the echo() procedure. Therefore, P₂ is suspended until P₁ comes out of the echo procedure.
 - Later on, process P₁ resumes and finishes the execution of echo procedure. The output displayed is the character m.
 - When P₁ comes out of echo() procedure, P₂ becomes active. Now process P₂ can successfully call the echo() procedure.

Therefore, if is necessary to use shared memory at a time.

3.2 Interprocess Communication

a Explain the inter-process communication in brief

MU - May 2016, 5 Mar

→ Ration and communication are two important elements should be satisfied when processes communicate with each other. Synchronization passes is required to achieve the mutual exclusion. Independent processes do not communicate with each other but cooperating processes may need to exchange information. Cooperative processes communicates through shared memory by message passing.

Cooperating processes require an Interprocess Communication (IPC) mechanism that will allow them to exchange data and information. There are fundamental models of interprocess communication.

Two fundamental models of Interprocess communication

- 1. Shared memory
- 2. Message passing

Fig. C3.1 : Fundamental models of interprocess communication

1. Shared memory

In the shared-memory model, a region of memory is shared by cooperating processes is established. Processes can then exchange information by reading and writing data to the shared region.

In the message passing model, communication place by means of messages exchanged between cooperating processes.

1 Message Passing

Message passing provides both functions. Message passing has the further benefit that it lends itself to implementation in distributed systems as well as shared-memory multiprocessor and uniprocessor systems.

- Following are the two primitives used in message passing :
 - o **send (destination, message)**
 - o **receive (source, message)**
- This is the minimum two operations required for processes to send and receive the messages. A process sends data in the form of a message to another process indicated by a destination. A process receives data by executing the receive primitive, indicating the source and the message.
- Communication by sending and receiving messages require synchronization. The receiver cannot receive a message until it has been sent by another process.
- The sending process is blocked until the message is received, or it is not after the send primitive executed by process. Similarly, when a process issues a receive primitive, **there are two possibilities :**
 - o Previously sent message is received and execution continues.
 - o If there is no waiting message, then either (a) the process is blocked until a message arrives, or (b) the process continues to execute, abandoning the attempt to receive.
- Thus, both the sender and receiver can be blocking or nonblocking. Three combinations are common, **although any particular system will usually have only** one or two combinations implemented :
 - o **The blocking send and blocking receive :** Until the message is handed over, the sender and receiver both gets blocked.
 - o **The nonblocking send and blocking receive :** After sending the message, the sender continues its work but receiver remains blocked until message is arrived to it. This permits a process to send one or more messages to a multiple destinations as quickly as possible. Here receiver is in need of message so that it can resume the execution . So it gets blocked until message arrives.
 - o **Nonblocking send, nonblocking receive :** Neither party is required to wait.

The nonblocking send and nonblocking receive:

 - o Both sender **and receiver** »'« " **wait** " will continue the work.

I Message passing system should give guarantee that messages will be correctly received by receiver. Receiver sends acknowledgement to sender after receiving the message, if acknowledgement not received in defined time then sender resend the message. It also offers authentication service.

3.3 Race Condition

Q- What Is race condition? Explain with example.

A race condition takes place when more than one process write data items so that the final result depends on the order of execution of instructions in the multiple processes. Consider two processes, P_1 and P_2 , which share the global variable "x". At the time of execution, process P_1 updates x to the value 1 and at the time of execution of another process, P_2 updates "x" to the value 2. Thus, the two tasks are in a race to write variable "x". In this example the process that modifies x value lastly decides the final value of "x".

Role of Operating System

1. The operating system should be capable of keeping track of the different processes.
2. Allocation and reclamation of various resources should be done by operating system.
3. Each process's data and physical resources should be protected by operating system against unintentional interfering by other processes.
4. The operations performed by the process and the output that it generates should not depend on the speed of execution compare to the speed of other concurrent processes.

Definition of process interaction

Process Interaction can be defined as :

- The processes not aware of each other
- The processes in a straight way or indirectly aware of each other

The conflict occurs between the concurrently executing processes. **if these processes compete for the use of the same resource.** More than one process may require the same resource while they are executing. The existence of the particular process is not known to other process. The competing processes do not pass the information to each other.

Engineering System (MU Sl- 4) [IT]
Syllabus Topic : The Critical Section Problem

3.4> The Critical Section Problem 161,
→ (Dec. 14, June 15)

Q.

ai section problem. nnr wri

MU - Dec. 2

Q.

Explain what is meant by critical section problem with its different solutions.

MU - June 15 10 Marks

As discussed previously, it is necessary to prevent more than one process from reading and writing the shared data at the same time. The portion of the program where the shared memory is accessed is referred to as the *Critical Section*. In order to avoid race conditions and be able to recognize them, the code in *Critical Sections* in each thread. The type of the code that comprises *Critical Section* are as follows:

- These codes reference one or more variables in a "read-update-write" way. At the same time, any of those variables may be changed by another thread.
- These codes modify one or more variables that can be referenced in "read-update-write" fashion by another thread.
- Any part of data structure used by the codes can be modified by another thread.
- Codes modify any part of a data structure that is currently in use by another thread.

When one process is currently executing shared modifiable data in its critical section, no other process is permitted to execute in its critical section. Hence, the execution of critical sections by the processes is mutually exclusive in time.

In a code called as critical section, only one process executes at a time. In a critical section problem an algorithm needs to be designed which permits at most one process into the critical section at a time, with no deadlock.

In a critical section's problem solution, must obey mutual exclusion, progress and bounded waiting. Any process directly cannot enter in critical section. First process has to obtain permission for its entry in its critical section. The segment of code which implements this is called the entry section. When process comes out of the critical section after completing its execution there it has to execute the rest of the code is the remainder section.

Process Control
Any solution to the critical-section problem must satisfy the following three necessary conditions:

Three necessary conditions of critical-section problem

1. Mutual Exclusion
2. Progress
3. Bounded Waiting

Fig . C3-2 : Conditions of critical section problem

→ 1. Mutual exclusion

At a time, single process should be executing in the critical section (CS). If currently Process P_1 is executing in CS then other processes should not execute in CS. Or

→ 2. Progress

If process's CS is free, means it is not executing in CS. In this situation suppose some processes go into their critical sections for execution, for example P_2 and P_3 . Only processes that are not executing in their sections are allowed to take part to make this decision. This decision should be taken in some definite time and should not be delayed indefinitely.

→ 3. Bounded waiting

A limit should be set on the number of times processes are permitted to go into their critical sections after a process has made a request to go into its section and before that request is approved.

Semaphore and monitor are the solutions to mutual exclusion.

A synchronization variable taking positive integer values is called semaphore. Binary semaphore only takes values 0 or 1. Hardware does not supply the semaphore. Semaphore offer the solution to critical section problem. Semaphore variable takes value greater than 1 then it is called as counting semaphore. Like semaphore, a monitor also solves critical section problem. It is a component which contains one or more procedures, initialization sequence and local data. Following are the components of monitors:

- Shared data declaration
- Shared data initialization
- Operations on shared data
- Synchronization statement

3.5 Mutual Exclusion

→ (June 15, Nov. 15, May 16)

Q. What is mutual exclusion?

MU - June 2015- Nov 2015. 5 Rmarks

Q. What is mutual exclusion? Explain its significance.

MU - May 2016, 5 Marks

At the same time as one process executes the shared variable, all remaining processes wishing to accomplish so at the same instant should be kept waiting; when that process has over executing the shared variable, one of the processes waiting to perform so should be permitted to carry on. In this manner, each process executing the shared data (variables) keeps out all others from doing so at the same time. Only one process should be allowed to execute in critical section. This is called Mutual Exclusion.

The mutual exclusion is put into effect only if processes access shared changeable data. If during execution, the operations of the processes do not conflict with each other, they should be permitted to proceed in parallel.

< Requirement of Mutual Exclusion

- It is obligatory to enforce a mutual exclusion: If many processes want to execute in critical section then only one process should be allowed to execute in that critical **section among all processes that have critical sections for the same resource or shared object.**
- **If the process halts in its remainder section (other than critical section) then it is allowed to do so without interfering with other processes.**
- **If process is waiting to enter in critical section then it should not be delayed for an indefinite period: that is it should not lead to deadlock or starvation.**
- If the critical section is empty (any process not **executing** in critical section), then if any process that requests entry to its **critical section must be allowed to enter without delay.**
- **Assumptions about relative process speed** Is or number of processors are not made.
- Any process will not remain inside its **critical section** for indefinite time. It should stay there or a **finite time** only.

Mutual Exclusion Conditions

A race condition can be avoided if we **disallow to two processes in their critical section at the same time.**

Following 4 conditions should hold to obtain an excellent solution for the critical section problem (mutual exclusion):

1. Two processes should not be inside their critical sections at the same time.
2. Comparative speeds of processes or number of CPUs in the systems are not assumed.
3. Process outside its critical section should not block other processes.
4. Any process should not wait for longer amount of time arbitrarily to enter its critical section.

Syllabus Topic : Peterson's Solution

3.6 Peterson's Solution

- Q. Explain Peterson solution to mutual exclusion.
Q. Give software approach for mutual exclusion.

- Peterson's solution involves only two processes. It gives a good means of algorithmic explanation of work out the critical-section problem and exemplifies some of the difficulties concerned in designing software that addresses the requirements of mutual exclusion, progress, and bounded waiting.
- The two processes P_i and P_j perform alternate **execution between their critical sections and remainder sections** (remaining code of the process). Peterson's solution requires the two processes to share two data items :

```
int x;
boolean y[2];
```

- The variable x decides who will enter in its critical section, if $x=i$, then process P_i is permissible to execute in its critical section. The y array is used to specify whether process w ready to enter its critical section. For example, if $y[i]$ is true, this value indicates that P_i is ready to enter its critical section.
- In order to enter the critical section, process P_i first sets $y[i]$ equal to true and then sets variable x to the value J . P_i $x=j$ because, if the other (P_j) process desires to enter the critical section, it can do so. If both processes attempt to enter simultaneously, x will be set to both i and j at approximately the same time. Only one of these assignments will last; the other will occur but will be overwritten straight away.

do {

```
y[0] = TRUE;
x = i;
while (y[i] && x == j);
```

critical Berben

y[i] = FALSE;

remainder section

} while (TRUE*

The final value of x decides which of the two processes is permitted to enter its critical section first. We can prove that solution is correct by showing:

- Mutual exclusion condition is satisfied.
- The progress requirement is satisfied.
- The bounded-waiting requirement also achieved.

➤ Mutual exclusion

- Every P_i goes into its critical section if and only if either $y[j] = \text{false}$ or $x = i$. Also note that, if both P_i and P_j runs in critical sections simultaneously, then $y[0] = y[j] = \text{true}$. It implies that P_i and P_j could not have successfully executed their while statements at about the same time, since the value of variable x can be either i or j but cannot be both.
- Hence, one of the processes say, P_i must have successfully executed the while statement, whereas P_j had to execute at least one additional statement (" $x = j$ ").
- However, at that time, $y[j] = \text{true}$ and $x = j$, and this condition will continue as long as P_i is in its critical section. Therefore we can conclude that mutual exclusion is preserved.

➤ Progress and bounded waiting

- Process P_i can be prohibited from entering the critical section only if it gets stuck in the while loop with the condition $y[j] = \text{true}$ and $x = j$; this is the only possible.

if $y[j] = \text{false}$, then P_i can go into its critical section. If $y[j] = \text{true}$ and $x = j$, then either variable x is set to j or $y[j]$ is set to true . In both cases, the process P_i will not enter the critical section.

Depending on the value of x either P_i or P_j will enter in the critical section.

If x is then true , it will enter the critical section.

After entering its critical section, it will release M so that P_j can enter its critical section.

If P_j resets $y[j]$ to true , it must also release M .

If P_j does not change the value of A while executing the while statement.

Critical section (progress) after at most (bounded waiting).

Syllabus Topic = Synchronization

3.7 synchronization Hardware

Hardware ut to mutual exclusion

- The simple hardware instructions many systems can be used successfully in solving critical section problem. Hardware features any programming task easier and get better efficiency.
- If we restrict the interrupts from taking shared variable was being modified, in environment the critical section problem is solved.
- Many systems have special instructions called Set Lock called as TSL instruction. It is having "TSL ACC, X". In this instruction ACC is accumulator register and X is symbolic name of memory location holds character to be used as flag.
- TSL is indivisible instruction means that it cannot be interrupted in between. After instruction gets following action take place.
 - o Content of X is copied to ACC
 - o Content of X becomes N

During execution of above two steps, TSL instruction is not interrupted. If we assume value of X=N initially, meaning of N is critical region is not free and F is free. Consider the following steps to enter in critical region.

1. TSL ACC, X
(Content of X is copied to ACC)
2. CMP ACC, "F"
(See if critical region is free)
3. BE 1
(If critical region not free then back.)
4. Return
(Return to caller and skip)

Following are the steps for exiting from critical region :

1. MOV X, "F* (F gets copied to X)
2. Return (Return to caller)

Consider the two processes P_i and P_j .

- Assume initially scheduler schedules P_i and $X = "F*$. So step 1 makes $ACC = F$ and $X = N$. Here P_i after executing step 4, prepare to enter CR (critical region).
- P_j gets scheduled due to context switch before P_i enters the CR. P_j executes step 1 and $X = N$. In step 2, P_j fails the comparison and loops back. It cannot execute step 3. So P_j cannot enter in CR as P_i is already in its CR.
- Consider P_i is again scheduled and executing in its CR. Even though context switch occurs and P_j gets scheduled at this moment it cannot enter in CR as ACC and X values are N (busy waiting).
- After completion of execution in CR, P_i executes step 5 to exit the CR and $X = **F$.
- If P_j is scheduled after this, it executes step 1 and ACC becomes F as X was F in above step. Now due to step 1 by P_j $ACC = "F"$ and $X = "F"$.
- Step 2 is executed by P_j and succeeds and it enters the CR.

This solution satisfies all conditions mutual exclusion, progress and bounded waiting. Since special hardware is required cannot be generalized to all machines. Also due to busy waiting, it is not efficient solution.

Syllabus Topic : Semaphores

3.8 Semaphores

Q, What is semaphore?

- E.W. Dijkstra (1965) abstracted the key notion of mutual exclusion in his concepts of semaphores.
- The solutions of the critical section problem represented in the section are not easy to generalize to more complex problems.
- To avoid this complicatedness, we can use a synchronization tool called a semaphore. A semaphore S is an integer variable that, apart from initialization, is accessed only through two standard atomic operations: wait and signal. These operations were firstly termed P (for wait) and V (for signal).

Definition

- A semaphore S is integer variable whose value can be accessed and changed only by two operations wait (P or sleep or down) and Signal (V or wakeup or up). Wait and Signal are atomic operations.
- Binary Semaphores does not assume all the integer values. It assumes only two values 0 and 1. On the contrary, counting semaphores (general semaphores) can assume only nonnegative values.
- The wait operation on semaphores S , written as $wait(S)$ or $P(S)$, operates as follows :

$wait(S): IF S > 0$

THEN $S: = S - 1$

ELSE (wait on S)

The signal operation on semaphore S , written as $signal(S)$ or $V(S)$, operates as follows :

$signal(S): IF (one or more processes are waiting on S)$

THEN (let one of these processes proceed)

ELSE $S: = S + 1$

- The two operations, wait and signal are done as single indivisible atomic operation. It means, once a semaphore operation has initiated, no other process can access the semaphore until operation has finished. Mutual exclusion on the semaphore S is enforced within $wait(S)$ and $signal(S)$.
- If many processes attempt a $wait(S)$ at the same time, only one process will be permitted to proceed. The other processes will be kept waiting in queue. The implementation of wait and signal promises that processes will not undergo indefinite delay. Semaphores solve the lost-wakeup problem.

cr Disadvantages of Semaphore

1. Semaphores are fundamentally shared global variables.
2. From any location in a program it can be accessed during course of execution.
3. A lack of command over it or assurance of proper usage.
4. A lack of proper linking between the semaphore and the data to which the semaphore controls access.
5. They provide two functions, mutual exclusion and scheduling constraints

Mutex is the concept related to binary semaphore. A key difference between the two is that the process that locks the mutex (sets the value to zero) must be the one to unlock

consumer do not access modifiable shared section of me buffer simultaneously.

Initialization

- Set full buffer slots to 0.
i.e. semaphore Full = 0.
- Set empty buffer slots to N.
i.e. semaphore empty = N.
- For control access to critical section set mutex to 1
i.e. semaphore mutex = 1.

Producer ()

WHILE (true)

- produce-item ();
- P (empty);
- P (mutex);
- enter-item ()
- V (mutex)
- V (full);

Consumer ()

WHILE (true)

- P (fill)
- P (mutex);
- remove-item ();
- V (mutex);
- V (empty);
- conaume-Item (Item)

Busy Waiting

► (Nov. 15, Dec. 16)

Q. What do you mean by busy waiting? What is wrong with it?

MU - Nov. 2015, Dec. 2016, 5 Marks

Busy waiting is the limitation of semaphore definition.

If critical section is not free (process is already executing in it) then other process trying to enter in critical section **should loop continuously in entry code. This busy waiting** (continual looping) is the problem as far as multiprogramming environment is considered. In this case a single CPU is shared among many processes. This busy waiting waste CPU cycles which would have been used by other processes effectively.

3.9.3 Readers/Writers Problem

Q. Explain reader-writer problem.

Process Coordination

- While defining the reader/writers problem, it is assumed that, many processes only read the file (readers) and many write to the file (writers). File is shared among a many number of processes. The **conditions that must be satisfied are as follows :**
 - o Simultaneously reading of the file is allowed to many readers.
 - o Writing to the file is allowed to only one writer at the same time.
 - o Readers are not allowed to read the file while writer are writing to the file,
- In this solution, the first reader accesses the file by performing a down operation on the semaphore file. Other readers only increment a counter, read count. When readers finish the reading, counter is decremented.
- When last one ends by performing an up on the semaphore, permitting a blocked writer, if there is one, to write. Suppose that while a reader is reading file, another reader comes along. Since having two readers at the same time is not a trouble, the second and later readers can also be allowed if they come down.
- After this assumes that a writer wants to perform write on the file. The writer cannot be allowed to write the file, since writers must have restricted access, so the writer is suspended. The writer will suspended until no reader is reading the file. If a new reader arrives continuously with very short interval and perform reading, the writer will never obtain the access of the file.
- To stop this circumstance, the program is written in a different way: when a reader comes and at the same time a writer is waiting, the reader is suspended instead of being allowed reading immediately.
- Now writer will wait for readers that were reading and about to finish but does not have to wait for readers that came along after it. The drawback of this solution is that it achieves less concurrency and thus lower performance. Following is the solution given.

```
typedef int semaphore;
semaphore mutex = 1 /* controls access to 'readcount' */;
semaphore file = 1; /* controls access to file */

int readcount = 0;

void reader(void)
{

```



```

while (TRUE) /* repeat evermore */
{
    down(&mu_ir_x); /* get exclusive */
    readcount = readcount + 1; /* one reader more now */
    if (rc == 1) down(&file); /* first reader */
    up(&ntu<<); /* release exclusive access */
    rd_fik(); /* read the file */
    downW(); /* exclusive to 'coun' */
    ree.ko.int = readreunt - 1; /* fewer */
    if (readout == 0) upf&iW; /* last reader */
    upfOrmutextV(); /* release exclusive access to 'readcount' */
    uee-datA-rradtV(); /* noncritical region */
}

void writer(void)
{
    while (TRUE) /* repeat evermore */
    {
        hink up data(); /* noncritical region */
        down(&file); /* get exclusive access */
        write_fileQ(); /* write the file */
        upl&file(); /* release exclusive access */
    }
}

```

3.9.4 Dining Philosopher Problem

→ (Dec. 14, Nov. 15)

- Q- Explain dining philosopher problem and solution to it [MU - Dec 2014, Nov. 2015. 10 Marks]
- Q. Explain Dining philosopher problem using semaphores.

- The problem is stated as : The work of five philosophers is thinking jinri patina a trtU« i *□ for them and all agreed that only food that contributed to their thinking efforts was spaghetti. Each philosopher requires two forks to eat spaghetti.
- The provision for eating** is a round table with a big serving pot of spaghetti, five plates, and five fork. One plate is given to one philosopher a-k: assigned place at the table. He has to use two forks on

either side of the plate. He takes and eats spaghetti-

Algorithm should be designed w, philosophers to eat. The algori ~~thm~~ must satisfy the exclusion (two philosophers should not fork at the same time) while avoidi ~~s~~ the starvation. This problem exemplifies basic deadlock and starvation.

The dining philosopher's problem is an example of problems related to coordination of resources, which may occur when it includes concurrent threads of execution in way, ~~h~~ Problem is a typical ~~w~~ca X approaches to synchronization.

Solution using semaphores

- In the solution using semaphores, each philosopher pick and choose the fork on the left side fork on the right side. Once the philosopher ends eating, the two forks are replaced on the table.
- This solution, unfortunately, leads to deadlock if the philosophers are hungry simultaneously, all pick up the fork on their left, and right side not be there. In this unseemly position, all philosophers starve.
- Additional five forks (a more clean solution), the problem or use of just one fork to eat the can be a solution. Following program shows solution.

```

/* program dining philosophers */
semaphore fork [5] = {1};
int i;
void philosopher (int i)
{
    while (true) {
        thiuk();
        wait (fork[i]);
        wait (fork [(i+1) mod 5]);
        eat();
        signalfork [(i+1) mod 5];
        signal(fork[i]);
    }
}
void main()

```

```
{
    parbe n (philosopher (0), philosopher (1),
    philosopher (2), philosopher (3),
    philosopher (4));
}
```

- In another strategy consider the helper who only allows four philosophers at a time into the dining room. With at most four seated philosophers, at least one philosopher will have access to two forks. Following is the solution using semaphore.

```
/* program dining philosopher* */
semaphore fork[5] = {1}; ■
semaphore room = {4};

int i;

void philosopher (ini i)
{
    white (true) {
        thinkQ;
        wait (room);
        wait (forkfij);
        wait (fork [(i+1) mod 5]);
        eat();
        signal (fork [(i+1) mod 5]);
        signal (forkfij);
        signal (room);
    }
}

void mainO
{
    parbegin (philosopher (0), philosopher (1),
    philosopher (2), philosopher (3),
    philosopher (4));
}
```

Syllabus Topic : Monitors

3.10 Monitors

Q. What is monitor? How it is used to achieve mutual exclusion? Explain.

Proc&ss Coordination

- If «maphnre s nTC uwd incorrectly, k efn pnjdue >*»nige™ that are hard to detect. These errors occur only if some particular execution sequences results and these sequences do not always occur. In order to handle such errors, researchers have developed high-level language constructs called as monitor.

- A monitor is a set of procedures, variables, and data structures that are all grouped together in a particular type of module or package. Processes may call the procedures in a monitor if required, but direct access to the monitor's internal data structures from procedures declared outside the monitor is restricted to the processes. Following is the example of the monitor.

monitor example

```
integer i ;
condition e ;
procedure producer ( );
end ;
procedure consumer ( );
end;
end monitor;
```

- Monitors can achieve the mutual exclusion: only one process can be active in a monitor at a time. As monitors are a programming language construct, the compiler manages the calls to monitor procedures differently from other procedure calls.
- Normally, when a process calls a monitor procedure, if any other process is currently executing within the monitor, it gets checked. If so, the calling process will be blocked until the other process has left the monitor. If no other process is using the monitor, the calling process may enter.

Characteristics of a monitor

- Only the monitor's procedures can access the local data variables. External procedures cannot access it.
- A process enters the monitor by calling one of its procedures.
- Only one process may be running in the monitor at a time; any other process that has called the monitor is blocked, waiting for the monitor to become available.

Conditional variables are contained within monitors. These conditional variables can be accessed only within the monitors and used to achieve the synchronization. Functions cwait() and csignal() carry out the operation on conditional variable and these are created as special data type in monitors.



1. **emdt (e)**: On execution of the cajW process gets poised. After access the monitor.
2. **csignal (0 > B ? e)**: P resumes. If blocked, Now ~~one~~ processes, were selected to resume execution. Only one ~~will be~~ processes, were the execution in case ~~any such~~ blocked. No action will be taken if there is no such process.
- Monitor wait and signal operations are dissimilar from those for the semaphore. If a process in a signals and no task is waiting on the condition variable, the signal is lost.
- When one process is executing in monitor, processes that trying to enter the monitor join a queue processes blocked waiting for monitor availability.
- Within the monitor process may provisionally block itself on condition x by issuing cwait (x); After blocking, it then waits in a queue of processes to enter again the monitor when the condition alters, and resume execution at the point in its program following the cwait (x) call.
- If the execution of process going on in the monitor and during execution, it notices alter in condition variable x, it issues csignal (x), which signals the related condition queue that the condition has altered.
- It is possible for producer to put in characters to the buffer only by way of using procedure append within the monitor; the producer does not have straight access to buffer.
- This procedure initially makes sure the condition nor M to conclude if there is space existing in the buffer ~~u space~~; if available then the process executing the monitor is blocked on that condition.

Syllabus Topic : Atomic Transaction

3 J2_Atomic Transaction

Q. What is atomic transactions? Explain,

- Critical means These
- Concurrent execution of two critical sections gives the result equivalent to their sequential unknown order. For data base system
- Important along with its storage and retrieval.

3.11.1 System Model of Atomic Transaction

While processing the transaction its ~~atomic~~ be preserved although there is possible TX there in and updates data within file on disk.

So transaction is treated as sequence of read and write operations that can be terminated by commit operation. Commit operation indicates is completed and terminated its execution. Whereas abort operation indicates that terminated unsuccessfully due to some error. As aborted transaction modifies the data at termination, the state of the data must be state what was before transaction started.

This roll back of transaction must be ensured by system. To ensure the atomicity, it is necessary to devices that store the data. Transactions are stored from these devices.

Volatile storage

Examples are main memory or cache memory. If system crash, information on this storage is not survived.

Nonvolatile storage

Examples are disk and magnetic tapes. If system crashes then information on this storage is not survived. Storage is slower than volatile storage.

Stable storage

Replicate the information on other disk to survive.

3.11.2 Log-Based Recovery

- The atomicity is ensured by recording modifications to the data by transaction. This is done on stable storage. Write-ahead logging method to achieve it.
- In this method, a log is maintained by system storage and each log record described single of by transaction write. Following are the fields of record.

o **Transaction name** : This name of ~~the~~

o ^{011*} ~~do~~ write operation is unique.

o ^{011*} ~~is~~ ~~Wn~~ ~~en~~ : This name of ~~data item~~ is also unique.

o ^{011*} ~~M~~ ~~value~~ ~~Th~~

- o New value : Value of data operation after write operation is carried out.
- Other log records are also there to record the significant events during processing of the transactions. These are for example, start of transaction, commit or abort of transactions. Using these logs the data is easily recovered in case of failures. Recovery algorithm uses following two procedures:
 - o **undo(Ti)** : The data items updated by transaction T_i is restored to old value by $\text{undo}(T_i)$ procedure.
 - o **rtndoTi** : This procedure set the value of data items to new value. These data items are updated by transaction T_i .

3.11.3 Checkpoints

Searching of log is carried out in case of failure. It is necessary to search the log to know the transactions to be redone and undone after failure. This searching is time consuming and modifications on updated data will cause recovery to get longer. Checkpointing is the solution over this problem. Along with maintaining the write-ahead log, system time to time takes checkpoints which involve the following sequence of actions to take place :

- Write all log records at present residing in main memory onto stable storage.
- Write all modified data residing in volatile storage to the stable storage.
- Write a log record <checkpoint> onto stable storage.

As <checkpoint> record is present in the log, it permits the system to make its recovery procedure more efficient.

3.11.4 Concurrent Atomic Transactions

Each transaction is atomic. There can be many transactions executing concurrently. This concurrent execution is equivalent to serial executions of the transactions in any of the order. This is called as serializability property. This property can be maintained by just executing each transaction within a critical section.

3.11.4(A) Serializability

Q. Write note on Serializability

- Let X and Y are two data items that can be read and written by two transactions TO and $T1$. Both transactions executes atomically in - the . order followed by $T1$. This sequence of execution is called as schedule. In serial schedule each transaction is

executed atomically. A serial schedule contains a sequence of instructions from different transactions in which the instructions of a particular transaction appear together.

If two transactions overlap their execution then the resultant schedule is no longer serial. A nonserial schedule does not essentially mean an incorrect execution. The Fig. 3.11.1 shows the serial schedule in which TO is followed by $T1$ The Fig. 3.11.2 shows the concurrent serializable schedule.

TO	T1
READ(X)	
WRIT(X)	
READ(Y)	
WRITE(Y)	
	READ(X)
	WRITE(X)
	READ(Y)
	WRITE(Y)

Fig. 3.11.1

Fig. 3.11.1 shows that serial schedule in which TO is followed by $T1$

TO	T1
READ(X)	
WRITE(X)	
	READ(X)
	WRITE(X)
READ(Y)	
WRfTE(Y)	
	READ(Y)
	WRITE(Y)

Fig. 3.11.2: A concurrent serializable schedule

Let OP_0 and OP_1 are the two operations from transactions TO and $T1$ respectively. If OP_0 and OP_1 accesses the same data item and at least one of the operation is write then OP_0 and OP_1 conflicts.

3.11.4(B) Locking Protocol

Q. Explain two-phase locking protocol.

- A lock is associated with each data item to ensure the serializability. In this case it is necessary by each transaction to follow the locking protocol that defines how locks are acquired and released. Following two models can be used.
 - o **Shared.** If a transaction T_i has acquired a shared mode lock on data item A , then T_i can read this item but cannot write A .

- Exclusive. If a transaction T_j has acquired an exclusive-mode lock on data item A, then T_i can both read and write A. *also ensures the lock and unlock requests in transaction should issue in two phases: growing phase, in this phase, a transaction may acquire locks but may not release any locks. Shrinking phase, in this phase, a transaction may release locks but may not acquire any*
- Two-phase locking protocol requires that, each transaction should issue lock and unlock requests in two phases: growing phase, in this phase, a transaction may acquire locks but may not release any locks. Shrinking phase, in this phase, a transaction may release locks but may not acquire any
- At the start a transaction is in the growing phase. The transaction acquires locks as required. After releasing a lock, it enters the shrinking phase, and no more lock requests can be issued.
- The conflict serializability is ensured by two-phase locking protocol but does not ensure freedom from deadlock. Beside this, it may happen that, for group of transactions, there are conflict-serializable schedules that cannot be obtained by using the two-phase locking protocol. To get better performance over two-phase locking, we must either need added information about the transactions or to impose some structure or ordering on the set of data.

3.11.4(C) Timestamp-Based Protocols

- In above two-phase locking protocol, the order followed by pairs of conflicting transactions is decided at execution time. Other method selects the order in advance to determine serializability order.

- For this time-stamp ordering scheme is used. System assigns a unique fixed timestamp to each transaction before it starts the execution. If TO transaction has

$T_{S(T0)}$ and later T_1 transaction enter in system having timestamp $TS(T_1) < TS(T_0)$. Following two methods are used to implement this approach.

- The value of the system clock can be used as the timestamp. In this method, the timestamp of the transaction is equal to the value of the clock when the transaction arrives in the system. If the transactions occur on separate processors, not sharing a common clock, the method will not work.

A logical counter is used as timestamp. In this method, the timestamp of transaction is the value of the counter when the transaction arrives in the system. The counter increases after a new timestamp is assigned.

Serializability order is determined by timestamp.

$TM_Z \gg$ $BTS < TO < TS < T$ \rightarrow to system produces equivalent schedule to serial schedule because BTS appears before T_1 .

Syllabus Topic : Deadlock

3.12 Deadlocks

► (June 15, Nov. 15, May 16, Dec. 2015, June 2016, Nov. 2016, May 2017, Dec. 2016, 2017)

Q. What is deadlock?

MU- June 2015, Nov. 2015, May 2016, Dec. 2016, 2017

Deadlock occurs in system when process waiting resources never gets resources as these are held by processes. So in this scenario, several processes change state from waiting to other state.

Syllabus Topic : System Model

3.12.1 System Model of Deadlocks

We know that processes needs different resources in order to complete the execution. So in a multiprogramming environment, many processes may compete for a large number of resources. In a system, resources are finite. With finite number of resources, it is not possible to satisfy the resource request of all processes.

When a process requests a resource and if the resource is not available at that time, the process enters a waiting state. In a multiprogramming environment it may happen with many processes. There is chance that

Waiting Processes will remain in same state and will never again change state.

It is because the resources have been requested are held by other waiting processes. When such type of situation occurs then it is called as deadlock.

If there are 5 same type of resources in the system then say that there are 5 instances of

- Suppose the process request the instance of resource. If any instance of the resource type is allocated to it, then request will be fulfilled. If it is not possible to fulfill the request, then the instances are not identical, and the resource type classes have not been defined appropriately.
- As a rule, first process must request a resource and then use it. After the use of this resource, it should release it. A process may request any number of resources which are necessary to complete its execution.

Following are the sequence when a process may utilize a resource :

- 1, **Request** : If it is not possible to fulfill the request instantly, then the requesting process must wait until it can get the resource.
- 2, **Use** : The process use the resource to accomplish the task.
- 3, **Release** : The process free the resource after use

Syllabus Topic : Deadlock Characterization

3.13 Deadlock Characterization

In deadlock, processes never finish execution as resources are held up by other processes. Due to this unavailability new processes are prevented from starting the execution. Following are the features that characterize the deadlock.

3.13.1 Conditions

In deadlock, processes never finish execution as resources are held up by other processes. Due to this unavailability new processes are prevented from starting the execution.

Necessary Conditions

→ (June 15, Nov. 15, May 16, Dec. 16)

- a. **Eti.--- necessary and sufficient condition**
deadlock to occur. **yMtpTjune 2015, Nov. 2015.
Fay 2016 Dec. 2016 g. gJVjjgkS,**

If the following four conditions hold simultaneously in a system then a deadlock situation can occur in the system. These are shown in Fig. C3.3.

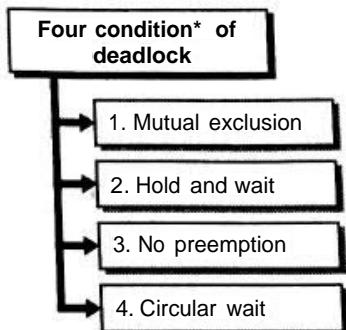


Fig. C33 : Necessary conditions of deadlock

♦ 1. Mutual exclusion

This condition ensure that, resource should be used by single process at a time and it remains with using process in non-shareable mode. If same resource is needed by another process then the requesting process must be postponed to use this resource. Requesting process will get the resource after released by holding process,

♦ 2. Hold and wait

The process holds minimum one resource and waits to obtain remaining needed resources that are at present being held by other processes.

→ 3. No preemption

Preemption of the resource will not be done; Process holding the resource, releases it on its own after the completion of its chosen task,

→ 4. Circular wait

In this condition, a collection of waiting processes say (P_0, P_1, \dots, P_n) exist 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 .

3.13.2 Resource Allocation Graphs

→ (Dec. 16)

- q write note on : Resource Allocation Graph. **MU - Dec. 2016. 5 Marks**

- A directed graph called a system resource-allocation graph can explain the deadlock precisely.
- In this graph the set of vertices V is partitioned into two different types of nodes. These are,

The set consisting of all the active processes in the system.

say $P = \{P_1, P_2, \dots, P_n\}$ is the set of all resource types in the system.

The set consisting of all the resources in the system.

say $R = \{R_1, R_2, \dots, R_m\}$ is the set of all resource types in the system.

In this graph if there is a path from process P to resource R and from resource R to process P then it means that an instance of resource type R has been allocated to process P .

- So, edge from process P to resource R and edge from resource R to process P is an edge. Circle nodes in graph represent processes and square represents resources.
- If resource allocation graph contains cycle then deadlock may exist in the system. If it does not contain cycle, then processes are not deadlocked.

If system contains single instance of all the resources

- then all processes involved in cycle is deadlocked. So in this case, a cycle in resource allocation graph is both necessary and sufficient condition for existence of deadlock. Dot inside square represents single instance of that resource.

- Cycle in following graph shows the deadlock situation when resource instances are single. P_1 has requested the resource held by P_2 and P_2 request the resource held by P_1 . P_3 has requested resource R_3 .

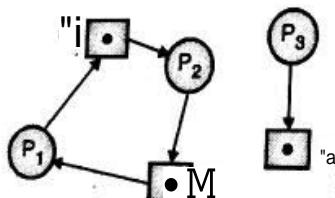


Fig. 113.1 : Resource Allocation Graph

- If more than one instances of each resource type is present in resource allocation graph then a cycle indicates that a deadlock has occurred. In this case a cycle is a necessary condition for the existence of deadlock.

- The reason behind cycle not being the sufficient condition is that, the process can release the held resource although cycle exists. This released resource then can be obtained by requesting process.

Conserve cycle can be broken. Hence, if present resource allocation graph then system may or may not be in a deadlocked state.

Syllabus Topic : Deadlock Prevention

3.14 deadlock prevention

→ (June 15, Nov. 2015)

Q. Explain deadlock prevention.

MU - June 2015; 2 Marks

- Q. What is difference between deadlock avoidance and deadlock prevention?

MU - Nov. 2015. 10 Marks

The described above, for a deadlock to occur, the four necessary-conditions: mutual exclusion, preemption and circular wait simultaneously in the system. By ensuring that at least one these conditions cannot hold, we can prevent occurrence of a deadlock. So deadlock prevention is achieved by denying the holding of at least one of, following four conditions.

1. Mutual exclusion
2. Hold and wait
3. No preemption
4. Circular wait

3.

The approach is explained as below.

1. Mutual exclusion

- Non-sharable resources cannot be shared by many processes simultaneously. If resources are sharable then mutual exclusion condition may hold for them.
- For example, it is not possible to use printer by many processes at the same time. Printer is non-sharable but read only file is sharable resource do not require mutually exclusive access.
- Therefore sharable resource such as read only cannot be involved in a deadlock. We can prevent deadlock by denying mutual exclusion as some resources are inherently non-sharable such as printer.

2. Hold and wait

- Here we have to see that hold and wait will not occur. The should be an assurance that process P_2 holds requests a resource only if it does.

- not hold any other resource. Following two protocols can be used.
- Allocate all the resources needed by process before it starts the execution. This protocol will keep the allocated resources unused. The reason is that, even though some of the resources are needed by process in the last stage of execution, it will hold it from start of execution.
- Second protocol allows the process to request the resources only when it has released the previously allocated resources. Again the disadvantage of first protocol appears here. Suppose the important resource for which many processes compete is needed only at the start and end of the execution by this process. If initially process release resource and at the end request for the same, then it may not possible to get the same resource immediately. So we must request all resources at the beginning.
- Apart from the limitation of the above two protocols, other limitation can be starvation. In this, processes will keep significant and popular resources with them forcing to wait other processes indefinitely which require these resources.

3. No Preemption

- We have to ensure that no preemption condition does not hold. Following two protocols can be used.
- If a process at present keep some resources with it and requests another resource that cannot be immediately allocated to it, then all resources at present being held should be preempted.
- The list of the resources for which process is waiting is maintained. The preempted resources from the process are included in this list. The process now restart the execution provided that it can get back its old preempted resources, in addition to the new ones that it is requesting.
- If the needed resources by the process are available, then allocate to the process requesting for it. On the other hand if they are not available, then check whether it is allocated to the process which is further waiting for additional resources. If it is the case then withdraw the resources from waiting process and allot them to requesting process.

- The process which requests the resources should wait in case If resources are neither available nor held by waiting process. In this situation, if other process make a request for the resource held by waiting process, then it should be preempted from waiting process.

4. Circular wait

- To prevent the occurring of circular-wait condition inflict a total ordering of all resource types. In this case each process should request the resources in an increasing order of assigned numbers to the resources.
- Consider the collection of resource types $R = \{R_1, R_2, \dots, R_n\}$. To this each resource type, a unique integer number is assigned. As resources are now recognized by their numbers, it is easy to compare them. It is also easy to know if one resource precedes the other in our ordering. Assignment of numbers to resources is a one-to-one function $p : R \rightarrow N$, where N is the set of natural numbers. In order to prevent the deadlock, following two protocols can be used.
 - Each process requests resources in an increasing order of assigned numbers to the resources.
 - Whenever process request instance of resource type R_j it should release any resource R_i such that $F(R_i) \geq F(R_j)$.

If these two protocols used, circular wait condition cannot hold. If process P_i waits for resource R_i which is held by process P_{i+1} we get $F(R_i) < F(R_{i+1})$ for all i . It implies that $F(R_0) < F(R_1) < F(R_2) < \dots < F(R_n) < F(R_0)$. By transitive relation, we get contradictory statement $F(R_0) < F(R_0)$. circular wait does not hold

" "s nTbusTopicTdeadlockAvoidance ~

3d5 Deadlock Avoidance

→ (June 15, Nov. 15)

- Q. Explain deadlock avoidance. MU • June 2015. 2 Marks
- Q. What is difference between deadlock avoidance and deadlock prevention? MU - Nov. 2015. 10 Marks

Resource avoidance requires that the system has some available up front. Initially every process needs the number of resources needed.

- An ItZn»»» method for avoiding deadlocks is require additional information about to be requested. If the sequence in which pr request the resource and will release it toow Jk.on.Hy, then we can decide for each reque. whether or not the process should wait.

To avoid the circular wait condition completely, a deadlock-avoidance algorithm inspects the resour allocation slate in dynamic way. Following are factors which define and describe resource allocation state.

 - o The number of available resources in the system.
 - o The number of allocated resources in the system.
 - o The maximum demands of the processes.

Dijkstra's Banker's algorithm is used for deadlock avoidance. The algorithm requires prior knowledge of number of resources each process needs.

After that, the OS act like a sincere small-town banker. OS will allocate the resources to processes only it has sufficient number of resources available to fulfill possible demand. This is considered as a safe state.

☞ Safe State

- If the system is able to allocate resources to each process up to its maximum need in some order and still avoid a deadlock then the system is in safe state.
 - If the safe sequence of all processes present in system the system is considered to be in safe state: Sequence $\langle P_1, P_2, \dots, P_n \rangle$ is considered as a safe sequence for the present allocation state if, for each process P_x , the resources which P_x can still request can be fulfilled by,
 - o The at present available resources in the system plus
 - o The resources held by all of the processes P_y 's, where $y < x$.
 - If process P_i cannot get all the needed resources as they are not available, it should wait until P_j complete the execution and frees the resources

If the system is a ~~deadlock~~ state, there can no possibility of deadlock. If the system is in an unsafe state, there is the

cr Example

Consider a system with 10 disk drives and 3 processes (P0, P1, and P2):

Process	Maximum needs	Current needs
PO	8	4
PI	4	2
P2	7	2

1.1710 disk drives are available in the line.

“fres 8 disk drives process P1 requi 'M

P2 may hold up to 7 disk drives. A4
ieTs holding 2 disk drives and process PO J

drives. (Thus in the system 2 disk drives are f

At time t_0 , the system is in safe state. 1 If
 $\langle p_1, P_0, P_2 \rangle$ satisfies safety condition, *!
immediately be allocated its disk drives 1 ii
them after completion. Now total 4 disk J 1
available. 1

- PO gets all the disk drives. After completion them and now total 8 disk drives are avail ! system. Finally P2 can be allocated all the ne J drives and returns them after completion. (The J will have now all 10 disk drives aJ Fig. 3.15.1 shows safe, unsafe and deadlock spaces.

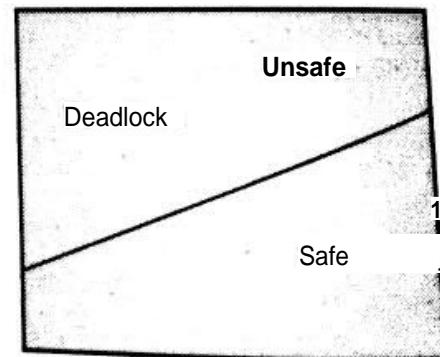


Fig 3.15.1 : Safe, unsafe, and deadlock states?

- At time t_h if any process requests the resource say disk drives in our example and if allocation other process could not be allocated **4** ~~resource type due~~ in unsafe state, then system t_0 unavailabilty, then system

3.1

5-1 Deadlock

Avoidance Algorithm#

→ (May ifc*

Following algorithms are used in the avoidance of deadlock.

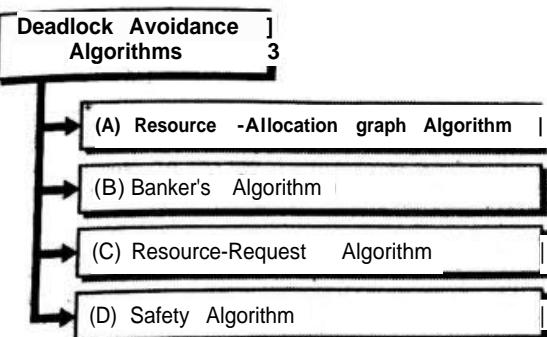


Fig. C3.4 : Deadlock avoidance algorithm

4 3.15.1(A) Resource-Allocation Graph Algorithm

- If system has resource allocation system with only one instance of each resource type, then only this algorithm is applicable.
- Future request edge in resource allocation graph is represented by claim edge (dotted edge).
- The claim edge is transformed to a request edge when a process requests a resource.
- The assignment edge is changed to a claim edge after a process free the resource.
- When requested resource is allocated to the process **then request edge is converted to assignment edge**. If assignment leads to cycle in resource allocation graph then system is in unsafe state. If not so then system is in safe state.
- After assignment of resource, if cycle detection **algorithm detects cycle, then system is in unsafe state**. The unsafe state in resource allocation graph is shown in Fig. 3.15.2.
- Cycle detection : $O(n^2)$

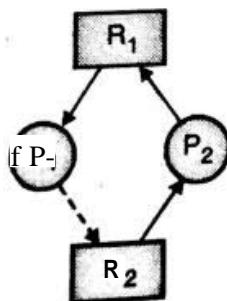


Fig. 3.15.2 ; Resource-allocation graph for deadlock avoidance

→ 3.15.1(B) Banker's Algorithm

→ (Dec. 16)

Q. Explain the banker's algorithm In detail.

MU - Dec. 2016, 10 Marks

- It is applicable to the resource allocation system with multiple instances of each resource type.
- It is less efficient than resource-allocation graph algorithm.
- Newly entered process should declare maximum number of instances of each resource type which it may require.
- The request should not be more than total number of resources in the system.
- System checks if allocation of requested resources will leave the system in safe state. If it will the requested resources are allocated.
- If system determines that resources cannot be allocated as it will go in unsafe state, the requesting process should wait until other process free the resources.

→ Data structures used In Banker's algorithm

→ (June 15)

Q. Explain data structures used in banker's algorithm.

MU - June 2015, 5 Marks

Following Data structures are used in Banker's algorithm

- **A[m]** : Array A of size m shows the number of available resources.
- **M[n][m]** : Two dimensional array M shows maximum requirement of the resources by each process. **M[i][j] = k** indicates process P_i can request at the most k instances of resource type R_j .
- **C[n][m]** : Two dimensional array C shows current allocation status of resources to each process. **C[i][j] = k** indicates process P_i allocated k instances of resource type R_j .
- **N[n][m]** : Two dimensional array N shows the remaining possible need of each process. i.e., $N[i][j] = M[i][j] - C[i][j]$. $N[i][j] = k$ indicates process P_i may need additional k instances of resource type R_j to accomplish the execution.

→ 3.15.1(C) Resource-Request Algorithm

- This algorithm decides whether requests of the resources can be safely granted so as to system will remain in safe state,

- Let $R[n][m]$ be the two dimensional array, describe the current outstanding requests of all processes. $R[i][j] \sim k$ indicates process P_i want k instances of resource type R_j to accomplish the execution.

1. If $R[i][j] \leq N[i][j]$, to step 2; otherwise, raise an error condition as process is requesting more resources than it needs.

2. If $R[i][j] > A[j]$ then the process should wait as needed resources are not available.

3. Otherwise, go to step 4

4. Allocate the requested resources to P_i . The state is altered as

$$A_{ijl} = A_{i,j-1}$$

$$C_{il[j]} = C_{il0} + R[i][j]$$

$$N_{ilU} = N_{ilj} - R_{ilj}$$

check to if the

- Once the resources are allocated, must wait system state is safe. If unsafe, the old resource-allocated states restored

→ 3.15.1(D) Safety Algorithm

state

- This algorithm find out whether system is safe or not.

initialized

1. U, w be an integer array of length m , $U[i] = 1$. Let F to array A (Available number of resource to false) (determines if process finished or not).

2. Find i such that both:

$F[i] = \text{false}$

$N[i] \leq w$

If D such i exists, go to step 4

3. $W = W + C(i);$

$F[i] = \text{true};$

Go to step 2

4. If $F[i] = \text{true}$ for all i , then the system is in a safe state, otherwise unsafe.

Run-time complexity : $O(m \times n^2)$

- To demonstrate the working of banker's algorithm, consider the following snapshot of the system at time t_0 . Processes are numbered P_0 to P_4 (total 5 processes).
- Resources are P, Q and R (total 3 resources). Resource type P has 12 instances, resource type Q has 9 instances, and resource type R has 6 instances.

Process	CfAIK****	Maximum requirement			N (Remaining need)
		P, Q, R	P, Q, R	P, Q, R	
P0	3. 0. 1	5, 2, 2		2, 2, 0	
P1	0, 2, 1	3, 2, 2		3, 0, 1	
P2	2, 0, <1	3, 2, 2	2	a-2	
P3	2, 1, 0	7, 3, 1	5, 2, 1		
P4	4, 0, 2	5, 0, 3	1, 0, 1		

,, currently with sequence $\langle p_1, p_2, p_3, p_4, p_0 \rangle$
The sys state- The sequence $\langle p_2, p_1, p_0, p_4, p_3 \rangle$ satisfy the safety criteria as shown.

satisfy the safety P 2 :

1. $p_0, (| 2, 2) \leq (6, 2)$ is true

$$N \leq A \quad (1, 6, 2) - (1, 2, 2) = (0, 4, 0)$$

$$A = A - C \quad (0, 4, 0) + (3, 2, 2) = (3, 6, 2)$$

$A = (3, 6, 2)$ available in the system after P_2 finishes.

2. process P_1 :

$N \leq A \quad (3, 0, 1) \leq (6, 2, 2)$ is true

$$A = A - N \quad (3, 6, 2) - (3, 0, 1) = (0, 6, 1)$$

$$A = A + C \quad (0, 6, 1) + (3, 2, 2) = (3, 8, 3)$$

available in the system after P_1 finishes.

3. process P_0 :

$$N \leq A \quad (2, 2, 0) \leq (3, 8, 3)$$

$$A = A - N \quad (3, 8, 3) - (2, 2, 0) = (1, 6, 3)$$

$$A = A + C \quad (1, 6, 3) + (5, 2, 1) = (6, 8, 4)$$

available in the system after P_0 finishes.

4. Process P_4 :

$$N \leq A \quad (1, 0, 1) \leq (6, 8, 4)$$

$$A = A - N \quad (6, 8, 4) - (1, 0, 1) = (5, 8, 3)$$

$$A = A + C \quad (5, 8, 3) + (5, 0, 3) = (10, 8, 6)$$

available in the system after P_0 finishes.

5. Process P_3 :

$$N \leq A \quad (5, 2, 1) \leq (10, 8, 6)$$

$$A = A - N \quad (10, 8, 6) - (5, 2, 1) = (5, 6, 5)$$

$$A = A + C \quad (5, 6, 5) + (7, 3, 1) = (12, 9, 6)$$

available in the system after P_0 finishes.

So the required resources can be allocated processes in the sequence $\langle P_2, P_1, P_0, P_4, P_3 \rangle$.

SHSE Suppose at time H . P? request one



msw®®⁶
tyP Q, !
, v, labl'

So

that reC

resotirc

wheth

exists. I

3.16

E®^r

Corn

instance of resource type P and two instances of resource type Q, so Request R2 - (1, 2, 0) and it is not greater than available resources A.

So step 4 of the resource request algorithm indicates that request can be allocated. After allocation of these resources, now safety algorithm is executed to determine whether system remains in safe state. If any such sequence exists, then system will be in safe state.

3.16 Solved Problems

Example 3.1.8.1 fa [*11 Dec. 2014] It-I^{II} Marks

Consider the following snapshot of a system :

	Allocation	Max	Available
	A B C, D	A B C, D	A B C, D
P0	0 0 1 2	0 0 1 2	1 5 2 0
P1	1 0 0 0	1 7 5 0	
P2	1 3 5 4	2 3 5 6	
P3	0 6 3 2	0 6 5 2	
P4	0 0 1 4	0 6 5 6	

With reference to banker's algorithm

- Find need matrix
- Is the system in a safe state?
- If a request from process P1 arrives for (0, 4, 2, 0), can the request be granted immediately ?

Solution :

- Need matrix is,

	Need			
	A	B	C	D
P0	0	0	0	0
P1	0	7	5	0
P2	1	0	0	2
P3	0	0	2	0
P4	0	6	4	2

- System is in safe state because many safe sequences exist such as (P0; P2; P1; P3; P4)- As available ≥ 1.5 . either process P0 or P2 can complete the execution. Once process P3 finishes, it releases resources which permit all other existing processes to execute.

- Yes, request can be granted immediately* After granting this, Available will become (1, 1, 0, 0). One of the safe sequence is (P0, P2, P3, P1, and P4). Allocation for P1 becomes (1, 4, 2, 0) and Need for P1 becomes (0, 3, 3, 0),

Example 3.16.2

Consider the following snapshot of the system.

Process	Allocation			Maximum Requirement		
	A	B	C	A	B	C
P0	0	0	1	1	1	1
P1	2	0	0	3	2	3
P2	1	3	2	4	3	1
P3	1	0	1	0	0	1
P4	0	0	1	3	2	1

Let ABC be the resources with instances of A is 7, B is 3 and C has 6 instances.

- What is the content of matrix need ?
- Is the system in a safe state? Prove it.

Solution :

- Need = Max - Allocation

A	B	C
1	1	0
1	2	3
3	0	1
-1	0	0
3	2	0

- $A = 7, B = 3$ and $C = 6$ are total instances of resources. By adding the Columns of Allocation matrix

We get, $A = 4, B = 3$ and $C = 5$,

Now Available number of resources = Total - Allocation

$$\bullet \text{ Resource A available} = 7 - 4 = 3$$

$$\text{Resource B available} = 3 - 3 = 0$$

$$\text{Resource C available} = 6 - 5 = 1$$

By using above available number of resources and applying bankers algorithm, sequence

$\langle P2, P3, P4, P0, P1 \rangle$ is safe sequence.

Hence, system is in safe state.

Example 3.16.3
Consider the system with 5 processes, p0 to p4 and resource type A, B, C. An rthe following questions.

Process	Allocation			Request	Available		
	A	B	C		A	B	C
P0	0	1	0	0	0	0	0
P1	2	0	0	2	0	0	0
P2	3	0	3	0	0	0	0
P3	2	1	1	1	0	0	0
P4	0	0	2	0	0	2	0

Let ABC be the resources with instances of A is 7, B is 3 and C has 6 instances.

- (i) Is the system in a safe state?
- (ii) If at time t_1 , P2 makes additional request for an instance of type C, Check whether system is in deadlock or Not at time t_1

Solution :

$$(i) \text{ Request } \leq \text{ Available } (0 \ 0 \ 0 \leq 0 \ 0 \ 0)$$

Hence request of P0 is satisfied. P0 free the resources. Hence request of P2 becomes (0 1 0)

to P0, Let the P2 request for resources.

$$\text{Request}(P0) \leq \text{Available} (0 \ 0 \ 0 \leq 0 \ 1 \ 0)$$

Hence request of P2 is satisfied. P2 free the resources (3 0 3) and Available becomes (3 1 3).

$$\text{Available} = \text{Available} + \text{Allocation} \{P2\}$$

In this way we can find that P3, P1 and finally P4 gets

all their requested resources and safe sequence exist $\langle P0, P2, P3, P1, P4 \rangle$. Therefore system is in safe state.

- (ii) If at time t_1 P2 makes additional request for an instance of type C then $\text{Request}[F0] \leq \text{Available} (0 \ 0 \ 1 \leq 0 \ 1 \ 0)$ is false and resource C cannot be granted to P2. System will be in deadlock state.

Example 3.164

It is proposed to use bankers algorithm for handling deadlock. Total number of resources available for allocation are 7, 7 and 10 respectively. The current resource allocation state is shown as below.

Process	Allocation			Max			Available
	R1	R2	R3	R1	R2	R3	
P1	2	2	3	3	6	b	7
P2	2	0	3	4	3	3	7
P3	1	2	4	3	4	4	7

19. the current allocation a safe state?

(i) the following request be fit? $\{1, 0, 1\}$

Solution :

Is it safe? $\{1, 0, 1\} \leq \{3, 6, b\}$

0. "d, x is

adding AH column

$$\text{Request} = \text{Max} - \text{Allocation}$$

Now available in system are (7, 7, 10)

The request of P2 satisfied as available in system $(7, 7, 10)$ and then of P3 is satisfied, the system is in safe state.

(ii) Request of P1 is satisfied as $\{1, 1, 0\} \leq \{7, 7, 10\}$

Example 3.165

Consider following snapshot of the system

Using Banker's algorithm answer the following questions

- (i) How many resources of type A B C D are there?
- (ii) What are the content of need matrix?
- (iii) Find if system is in safe state? If it is, find sequence.

Process	Max				Allocation				Available		
	A	B	C	D	A	B	C	D	A	B	C
P0	6	0	1	2	4	0	0	1	3	2	1
P1	1	7	5	0	1	1	0	0			
P2	2	3	5	6	1	2	5	4			
P3	1	6	5	3	0	6	3	3			
P4	1	6	5	6	0	2	1	2			

Solution :

$$\langle A-9; B-13; C-10; D-11 \rangle$$

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

Need matrix is



	A	B	C	D
PO	2	0	1	1
P1	0	6	5	0
P2	1	1	0	2
P3	1	0	2	0
P4	1	4	4	4

- (iii) The system is in a safe state as the processes can be finished in the sequence PO, P2, P4, P1 and P3,

Example 3.16.6 MU - June 2015. Nov. 2015. 10 Marks

Consider the following snapshot of the system

	Allocation	Max	Available
	A B C	A B C	A B C
P0	0 1 0	7 5 3	8 3 2
P1	2 0 0	3 2 2	•
P2	3 0 0	• 0 2	•
P3	2 1 1	2 2 2	•
P4	0 0 2	4 3 3	•

Answer following questions using banker's algorithm

- What is the content of need matrix?
- Is the system in a safe state?
- If a request from process P1 arrives for (1, 0, 2), can the request be granted immediately?

Solution :

$$(a) \text{ Need} = \text{Max} - \text{Allocation}$$

	A	B	C
PO	7	4	3
P1	1	2	2
P2	6	0	0
P3	0	1	1
P4	4	3	1

- System is in safe state because safe sequence exist and it is (P1; P3; P4; P2; PO)-
- Yes, request can be granted immediately The reason is Request \leq Available that is (1, 0,) < (1, 0, 2)

Example 3.16.7 MU - Nov. 2015. 10 Marks

Consider the snap of the system.

Process	Allocation A B C D	Max				Available A B C D
		a	b	c	d	
P0	0 2 1 2	0	3	2	2	2 5 3 2
P1	1 1 0 2	2	7	5	2	•
P2	2 2 5 4	2	3	7	6	•
P3	0 3 1 2	1	6	4	2	•
P4	2 4 1 4	J	3	6	5	8

Answer the following questions using bankafs algorithm.

- What is the content of Matrix need?
- Is the system in safe state?
- If a request from process P1 arrives for (1,3,2,1) can the request be granted immediately?

Solution :

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

So content of Need matrix is

	A	B	C	D
P0	0	I	1	0
P1	I	6	5	0
P2	0	I	2	2
P3	1	3	3	0
P4	1	2	4	4

- Yes, System is safe and safe sequence (PO, P2, P1, P4, P3) exist.

- As (3, 3, 4, 8) remains available in system, P1 can be granted (1,3,2,1) immediately

Sylabus Top7c : Deadlock Detection

3.17 Deadlock Detection

→ (June 15)

Q. Explain deadlock detection.

MU - June 2015. 2 Marks

If the system does not make use of deadlock prevention or deadlock detection algorithm, then system will observe the **deadlock situation**. In this situation, it is necessary that system should have

3.18 Deadlock Recovery

Q. Explain deadlock recovery in detail.

- It needs to recover from deadlock when it is detected. Several options exist after detection algorithm detects that a deadlock exists in the system. One possibility is to inform operator and let them decide how to deal with it manually.
- The second option is to allow the system to recover from the deadlock on its own (automatically). There are two solutions exist to break a deadlock. One solution is just to abort one or more processes so that circular wait will be broken. The second option includes preemption of some of the resources from one or more of the processes which are involved in deadlock.

3.18J Process Termination (Kill a process)

When deadlock occurs, the operating system decides to kill one or more processes. This will reclaim the resources acquired by that killed processes. The deadlock detection algorithm detects whether deadlock exist or not after killing the process. Following two methods are used for killing the process :

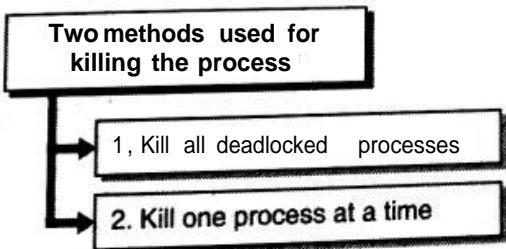


Fig. C3.5 : Two methods used for killing process

→ F 1. Kill all deadlocked processes

- This method is simple and effective to eliminate the deadlock and clearly will break the deadlock cycle, but at a huge cost.
- If all these killed processes have been performed die computation for longer period of time, then the partly computed result would be wasted. Recomputation will be done and is very costly.

→ 2. Kill one process at a time

In this method, one process is killed at a time and check whether detection algorithm is in ^{vo} _{to} deadlock is still exist or not in the system

- Invocation of the detection algorithm after killing each process is considerable overhead. We must re-run algorithm after each kill.
- Killing a process is not easy. We should kill only those processes whose killing will incur minimum cost.
- Different parameters that need to be consider to ensure the minimum cost are :
 - o Priority of the process.
 - o How much computation process has finished and how much more time does it need to finish the execution.
 - o The number of the resources process has used.
 - o The type of resources the process has used.
 - o The number of the resources the process will need to complete the execution.
 - o The number of processes needs to be killed.
 - o Whether the process is interactive or batch.

3.18.2 Resource Preemption

Incrementally preempt the resources from the process and reallocate resources to other processes until the circular wait is broken.

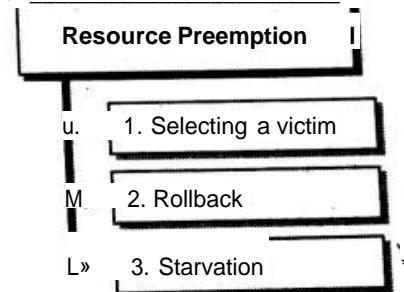


Fig. C3.6 : Resource preemption

→ 1. Selecting a victim

Which process and which resources should be selected for preemption to minimize the cost? The process which has finished almost all computation and only few amount of computation is remained should not be selected as a victim.

→ 2. Rollback

After preempting the resources from the particular process, it cannot continue the execution as needed resources are preempted. It is required to rollback the process to safe state and restarts the execution from that state. As safe state is difficult to define, total rollback will be the simple solution.

What do you mean by busy waiting? What is it? (Refer section 3.9.2) (5 Marks) (Nov. 2015)

3. Synchronization
It may W jesuittes-
times io p «P the

3.19 Exam pa an dRey! ?3
(University)

Synchronization

syHabutToP synchronization in brief. (May 2016)

Q. SIES 5 *

Q. Explain the inter-process communication in brief. (May 2016)

Q. Explain the inter-process communication in brief. (Refer section 3.2) (5 Marks)

Q. " " " with example

Q. " " " (Refer section 3.3)

Q. Syllabus <°Pc : The Critical Section Problem

a. Explain critical section problem. (Dec. 2014)

(Refer section 3.4) (5 Marks)

a. Explain critical section problem. With its different solutions. (Refer section 3.4) (10 Marks) (June 2015)

a. What is mutual exclusion?

Marks)

<»*-*"

3.9(5)

(June 2015, Nov. 2015)

Q. What is mutual exclusion ? Explain its significance.

(Refer section 3.5) (5 Marks)

(May 2016)

Q. Syllabus Topic : Peterson's Solution

Q. Explain Peterson solution to mutual exclusion. (Refer section 3.6)

Q. Give software approaches for mutual exclusion. (Refer section 3.6)

Q. Syllabus Topic : Synchronization Hardware

Q. Explain hardware solution to mutual exclusion. (Refer section 3.7)

Q. Syllabus Topic : Semaphores

Q. What is semaphore? (Refer section 3.9)

Q. Syllabus Topic : Classic Prnh tefl. of Synchronization 0 ems of

Q. Explain an algorithm for aoril ... er-consumer problem. (Refer section 3.9.1) (10 Marks) (Dec. 2016)

Q. Explain how Realtor with semaphores? (Refer section 3.9.2) (10 Marks)

Q.

What do you mean by busy waiting? What is it? (Refer section 3.9.2) (5 Marks) (Nov. 2015)

Q. Explain «ad riterproWem

Q. (Refer section 3.9.3)

Q. explain < 9 philosopher problem and Marfa) it. (

Q. Explain 0Wn9 philosopher Drnk semaphores. (Refer section 3.9.4) *1

Q. Syllabus Topic : Monitors

Q. YYtratis monitor? How it is used to ach ... exclusion? Explain. (Refer section 3.10) >

Q. Syllabus Topic : Atomic Transaction,

Q. What is atomic transactions? Explain. (Refer section 3.11)

Q. Write note on Serializability.

(Refer section 3.11.4(A))

Q. Explain two-phase locking protocol.

(Refer section 3.11.4(B))

Q. Syllabus Topic : Deadlocks

Q. What is deadlock ? (Refer section 3.12) (2 Marks) (June 2015, Nov. 2015, May 2016, Dec. 2015)

Q. Syllabus Topic : Deadlock Characterization

Q. Explain necessary and sufficient conditions for deadlock to occur. (Refer section 3.13.1) (2 Marks) (June 2015, Nov. 2015, May 2016, Dec. 2015)

Q. Write note on : Resource Allocation Graph.

(Refer section 3.13.2) (5 Marks) (Dec. 2015)

Q. Syllabus Topic : Deadlock Prevention

Q. Explain deadlock prevention. (Refer section 3.14) (2 Marks) (June 2015)

Q. What is difference between deadlock avoids and deadlock prevention?

(Refer section 3.14) (W Marks) (W*)

Q. II, b, U, T opic : O «k »k » «*! . «.

Q. EW «! a ~ »ad.,,,C..

(Refer section 3.14)

Q. What is difference between deadlock avoidance and deadlock prevention?
(Refer section 3.15) {1Q Marks} (Nov. 2015)

Q. Suggest techniques to avoid deadlock.
(Refer section 3.15.1) (1 Mark)
(May 2016, Dec. 2016)

Q. Explain the banker's algorithm in detail.
(Refer section 3.15.1(B)) (10 Marks) (Dec. 2016)

Q. Explain data structures used in banker's algorithm.
(Refer section 3.15.1(B)) (5 Marks) (June 2015)

Example 3.16.1 (10 Marks) **(Dec. 2014)**

Example 3.16.6 (10 Marks) (June 2015, Nov. 2015)

Example 3.16.7 (10 Marks) (Nov. 2015)

☛ Syllabus Topic : Deadlock Detection

Q. Explain deadlock detection.
(Refer section 3.17) (2 Marks) (June 2017)

Q. Write note on "Deadlock Detection Algorithms".
(Refer section 3.17))

☛ Syllabus Topic : Recovery from Deadlock

a. Explain deadlock recovery in detail
(Refer section 3.18) -----qq'

CHAPTER

4

Memory Management

Syllabus Topics

Memory Management strategies : Background, Swapping, Structure of the Page Table, Segmentation; Virtual Memory Write, Page Replacement, Allocation of Frames, Thrashing, Memory-Mapped Files, Allocating Kernel Memory, Other Considerations.

Syllabus* Topic : Memory Management Strategies

4.1 Memory Management Strategies

Background

Q. What are different requirements for memory management ? Explain.

- Memory is an important resource of the computer system that needs to be managed by the operating system.
- To execute the program, user needs to keep the program in main memory. The main memory is volatile.
- Therefore a user needs to store his program in some secondary storage which is non volatile.
- It is required to have process code, stack, heap (dynamically-allocated structures), and data (variables) to be in primary memory.
- Therefore memory must be allocated to every process. The management of main memory is required to support for multiprogramming.
- Many executables processes exist in main memory at any given time.
- Different processes in main memory have different address space.
- Memory manager is the module of the operating system responsible for managing memory.
- Programs after completion of execution move out of main memory or processes suspended for completion

can be swapped out on secondary storage to free in memory for other processes.

- New processes are required to be loaded in main memory.
- If available main memory is not sufficient to hold all processes, swapping between main memory and secondary memory is done.
- Memory managers move processes back and forth between memory and disk during execution.
- So it is required that operating system should have some strategy for the management of memory.

4.1.1 Monoprogramming

- In early computer systems and early microcomputer operating systems, mono-programming concept was used.
- Only one program was residing in main memory at given point of time.
- Operating system was residing in some portion of memory and other portion was fully devoted to single executing process.
- This system was simple to design but was supporting multiprogramming. In this approach relocation was needed as program could be loaded at the same location.

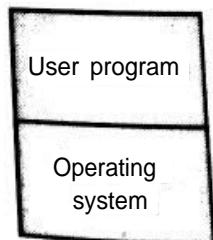


Fig. 4.1.1 : Monoprogramming



4.1.2 Multiprogramming

- Multiprogramming is required to support multiple processes simultaneously. Since multiple processes are resident in memory at the same time, it increases processor utilization if the processes are I/O bound, CPU does not remain idle.
- CPU utilization increases as it executes multiple processes on time sharing basis.
- Multiprogramming gives illusion of running multiple processes at once and provides users with interactive response to processes.

4.1.3 Dynamic Loading

- Dynamic loading ensure the better memory space utilization. Unless and until called, routine is not loaded in main memory in dynamic loading.
- All routines remain stored on disk in a relocatable load format.
- The main program is loaded into main memory and is executed.
- The advantage of dynamic loading is that, memory space is utilized only for the routines which are currently need to be executed and other unused routines reside on disk.
- If user designs the program which uses many library routines then operating system will itself support for dynamic loading. No special support require from operating system.

4.1.4 Overlays

- In the early days, as Large programs were too large to fit into their partitions, programmers created overlays to solve this problem.
- Rather than loading an entire program into memory at once, we can divide it into those modules that are always required and those that are required only at certain times during program execution.
- Overlays involve moving data and program segments in and out of memory which helps for reusing the area in main memory.
- Overlays were physical divisions of the program that ^{were} established by the programmer. To complete the execution, it is required that, entire program and data of ^a Process must be in physical memory.

- It is required to have process in main memory for execution and hence physical memory should accommodate the process.
- The process size can be larger than the amount of memory allocated to it.
- The overlays allows to have needed modules and data in main memory at any given time.
- The data and modules which are not required for execution at that time are swapped out of memory.
- Next time, the needed modules are loaded into space which is freed now by swapping the modules not required for execution.
- An overlay handler exists in another area of memory and is responsible for swapping overlay pages or overlay segments.
- Overlays are implemented by user; no special support needed from operating system, programming design of overlay structure is complex.
- Consider the example of two pass assembler. As we know that pass 2 is executed after pass 1, it is not necessary to keep both modules of pass 1 and pass 2 in memory simultaneously.
- During pass 1 symbol table is constructed and pass 2 generates machine-language code.
- Symbol table and common routines are required in both passes. Consider the memory requirement of pass 1 module, pass 2 module, symbol table and common routines.

Pass, 1 Module	100 K
Pass 2 Module	90 K
Symbol Table	30 K
Common Routines	30 K
Total	<u>250 K</u>

- Total 250 K memory is needed to run the assembler. If only 200 K memory is available then overlays can be defined as overlay A and overlay B.
- Overlay A contains pass 1 module, symbol table and common routines. Overlay B contains pass 2 module, symbol table and common routines.
- Overlay driver of 20 K is loaded in memory. Initially overlay A is loaded in memory which requires 160 K of memory.



- Once overlay A finishes execution, overlay B in place of overlay A. Overlay B requires IM K of memory.

4.1.5 Relocation

a) write short note on Relocation.

Consider the partitions of memory as follows.

- First partition: 100K (Operating system)
 - Second partition: 100K
 - Third partition: 20dK
 - Fourth partition: 300K
- Consider very first instruction "fci" file created by absolute address 100 within the binary
- UXI loads the program in first partition at address KJOK. Then instruction will find OS at absolute address 100.
 - If instead of first partition, loader loads the program in second partition, the instruction will jump to address 100K+100.
 - For third partition instruction will jump to address 200K+100 and so on. This problem is called as relocation problem.
 - The solution to this problem is to in fact modify the instructions as the loading of program is done in memory.
 - If programs are loaded in first partition then 100K is added to each address, programs loaded into second have 200K added to addresses, and so on.
 - To carry out relocation at the time of loading like this, the binary program must have linker included list or bitmap indicating which addresses are to be relocated and which are opcodes, constants, or other items not needed to relocate.

4.1.6 Logical and Physical Address Space

Q. With the help of examples, clearly differentiate between logical and physical address space.

- An address generated by CPU is called logical address or virtual address and an address generated by Memory Management Unit (MMU) is called physical address.
- If size of program written by user is 1000 K and loaded in main memory from address 5000 to 6000.

Physical address space will become 5000 to 6000 and logical address space is 0 to 1000. The collection of logical addresses JP* generated by program and a physical address J* correspond to physical addresses.

the 8 byte bmdirty¹ instructions "d" to tens. 1

When the same time or load address are same. If the same block occurs at run time then logical and physical address are different.

TM location register (base register) contains address.

IX) Physical address (5200) can be calculated.

Zws. Consider logical address 200 of program.

Physical address (5200) = logical address (offset) + content of base register (5000)

Moving from virtual to physical addresses is done by the memory management unit (MMU) at run time. The user program deals with logical addresses. The memory-mapping hardware converts logical address into physical addressed logical addresses.

Memory protection is done by using two registers, usually a base and a limit, as shown in Fig. 4.1.2. The base register holds the smallest legal physical memory address (in our example 5000). The limit register gives the size of the range of physical addresses. For example, if the base register holds 5000 and the limit register is 1000, then the program can legally access all addresses from 5000 through 6000.

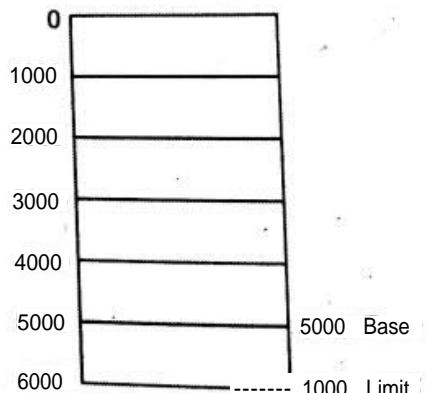


Fig. 4.1.2 : Memory Protection with base and limit registers

Syllabus Topic : Swapping

4.1.7 Swapping

- Systems in which it is possible to move entire Presses between disk and main memory during
- These are called swapping systems.

- In a time sharing operating system, system's memory is allocated to multiple processes.
- In order to have sufficient memory free, it will require to swap data between memory and secondary storage (disk).
- If we use round robin CPU scheduling algorithm then after expire of time quantum memory manager will swap out the finished process and swaps in the another process.
- It increases degree of multiprogramming ensures effective management of memory among many processes.
- If scheduling is priority based, then after arrival of high priority process low priority process is swapped out on disk.
- After completion of high priority process again a low priority process gets swapped in memory by memory manager.
- Swapping needs a backing store to store a swapped data. The backing store is usually a fast disk and having capacity enough to store copies of all memory images for all users.
- This backing store must offer direct access to these memory images.
- The operating system maintains a ready queue for the processes which are ready to execute.
- These processes have their memory images on the backing store or in memory.
- The CPU scheduler calls the dispatcher program, when it decides to run the process. If dispatcher does not find next process in the queue in main memory, it swaps in the desired process.
- If memory region is not free to bring the process from backing store, dispatcher swaps out current process from the memory.
- Then reloads registers and transfers control to the elected process. The context-switch time in such a swapping system is quite high.

Syllabus Topic : Chapter-11 Memory Allocation

4.2 Contiguous Memory Allocation

4.2.1 Multiprogramming with Fixed and Variable Partitions

- Q. Explain memory management with Fixed and Variable Partitions.
- Q. What is internal fragmentation and external fragmentation? Explain with example.

- In this scheme memory is divided into a number of fixed-sized partitions. Each partition may contain exactly one process.
- The number of partitions and its size would be defined by the system administrator when the system starts up. The degree of multiprogramming will be high if the number of partitions created is more.
- Always a selected process from the input queue is placed into the free partition. When the executing process finishes the execution and terminates, the partition becomes free for another process.
- The operating system maintains a table to keep track on free and occupied partitions by processes.
- Partitions can be of fixed size or of variable size. Fixed sized partitions are relatively simple to implement. The problems with the fixed size partitions are :
- If program size is larger than the partition size, program cannot be stored in partition.
- The second problem is **internal fragmentation**. If partition size is larger than program size, some space of the partition will be unused within that partition.

Fig. 4.2.1 shows fixed size partitions.

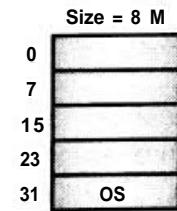


Fig. 4.2.1: Fixed size partitions

- As the name suggests, in a variable-sized partition, partitions of different sizes are available. Memory is **partitioned in partitions of different sizes**.
- The best-fit strategy is used to allocate the partitions to the processes. In best-fit, that partition is chosen to fit the process so that less amount of memory will be wasted.
- It means smallest partition big enough to accommodate the process is chosen.
- A queue is organized for each size of the partition where processes wait for that partition.
- Partitions are allocated to the processes with best-fit policy. Best fit policy ensures to allocate the partition **that is big enough for the process to reside**.
- Therefore variable partitions minimize internal fragmentation.



- However, such an allocation may reduce the performance as processes would queue up **in best fit** queue even though other partitions are free.
- For example, it may have many processes queued up in a queue of partition of size 8, and no processes are queued up in a queue for the partition size 16M.
- In both fixed and variable size partitions it may happen that, partition can remain empty even though processes are waiting in other partitions queue leading to external fragmentation. Fig. 4.2.2 Shows variable size partitions.

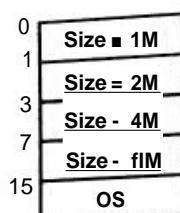


Fig. 4.2.2 : Variable size partitions

Difference between Internal fragmentation and external fragmentation

Q. Compare internal and external fragmentation.

Sr. No.	Parameter	Internal Fragmentation	External Fragmentation
1.	Partitioning technique	Internal fragmentation takes place in fixed partitioning technique	External fragmentation takes place in variable partitioning technique
2.	Definition	If partition size is larger than program size, some space of the partition will be unused within that partition called as internal fragmentation	in both fixed and variable size partitions it may happen that, partition can remain empty even though processes are waiting in other partitions queue leading to external fragmentation.

Sr. No.	parameter	Internal ¹ Fragmentation	External Fragmentation
3.	Drawback	Contiguous memory allocation reduce the internal fragmentation	Paging and segmentation technique is the solution over external fragmentation problem.

4.2.2 Dynamic Partition Technique

- Dynamic partition technique can reduce the internal and external fragmentations.

- In this technique basically a variable partitioning, variable number of partitions determined dynamically. **Initially, all memory is available for user processes, and is considered as one large block, of available memory, a hole.**
- When a process arrives and needs memory, it search for a hole large enough to fit for this process.
- This technique ensures the efficient use of main memory. There is no internal fragmentation.
- However, if the many small holes are scattered, it can't be allocated to one large size process. Consider the following example of Fig. 4.2.3.
- Initially 64 MB main memory is available. 8 MB is allocated to operating system as shown in Fig- 4.2.3(a). For user processes 56 MB memory is available and considered as one large single hole.
- Processes P1, P2 and P3 are loaded and to tiies processes sufficient memory space can be allocated, creating the hole of size 4 MB at the end of memory (Figs. 4.2.3(b), (c) and (d)).
- At some point of time none of the allocated processes are ready.
- Operating system swaps out P2 and in its place P4 of size 18 MB is swapped in Fig- 4.2.3(c).
- The hole of 2 MB is created. Again after some time since none of the processes in main memory is ready and only P2 is in ready suspend state, is available.
- To swap in P2, there is no sufficient space is available.
- So operating system swaps P1 out and swaps P2 back in as shown in Fig. 4.2.3(f).

Sh X plesbowthat,hereareman

small y scattered

These holes cannot be allocated to the size than these holes. This is fragmentation.

Process of larger
called external

Memory Management

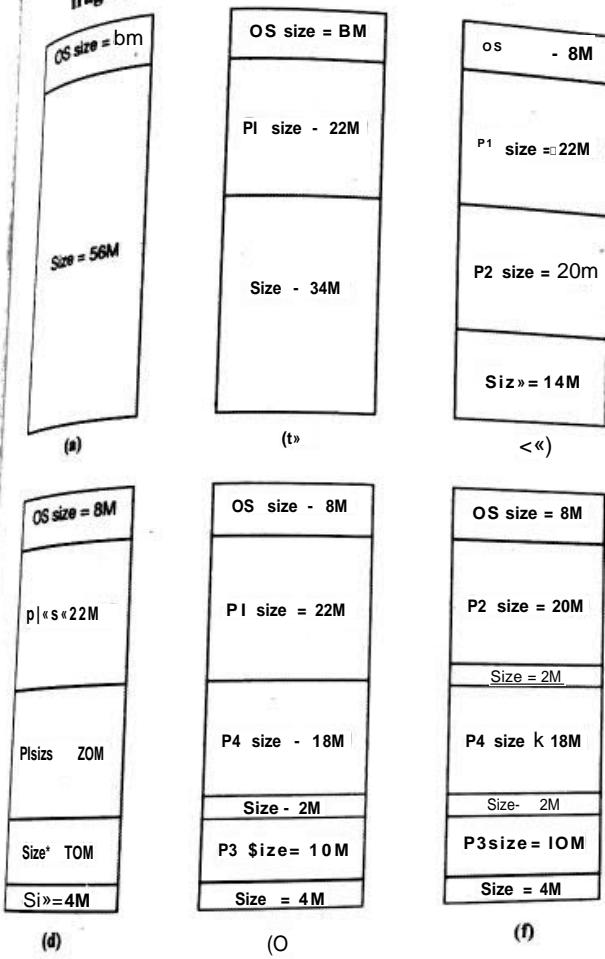


Fig.W : Creation of scattered holes in dynamic partition technique

2.3 Compaction

Q. What is compaction? Explain with example.

Solution to above discussed problem is compaction. Compaction is a technique to convert many small size holes into one large size continuous hole. Altered holes into one large size continuous hole. location of the various processes so as to accumulate all the free space in one place. Due to compaction, there is efficient use of process. Fig. 4.2.3 shows the memory after compaction. 4.2.4(a) shows the scattered small size holes. If processes P4 and P3 moved up, the end of memory as shown in Fig.

- If processes P2, P4 and P3 one 16 MB size hole of 8 MB is moved downward, we get area as shown in Fig. 4.2.4(c).
- However in this, 48 MB processes need to be moved compared to 28 MB shown in Fig. 4.2.4(b). Fig. 4.2.4(d) shows the 28 MB shifting of processes.
- Compaction should ensure the moving of processes having minimum total size. Compaction increases the degree of multiprogramming.
- It is because processes can be loaded in the memory which is available by combining many scattered holes. If relocation is done at compile time (static), compaction is not possible.
- It is possible only if relocation is dynamic and done at execution time.

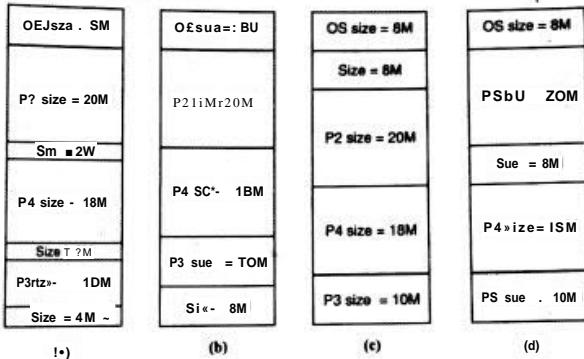


Fig. 4.2.4 : Compaction

4.2.4 Memory Allocation Strategies

-> (May 2016)

- Q. Discuss partition selection algorithm in brief. (MU ~ May 2016, 4 Mares)
- Q. Explain different memory allocation strategies.

If process requests the free hole, then question is which free hole should be allocated to this process. Following are the strategies used to allocate the free hole to the requesting process.

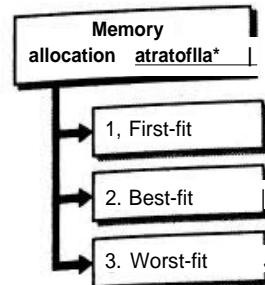


Fig. C4.1 = Memory allocation strategies

other pages are not

le to us and remains

do**1

in the same manner, we can say even

We have a large program available, the user only gasman set of instructions to execute at any time.

In fact, all these instructions which the processor needs to execute are within a small proximity of each other. It is like a page which contains all the statements which we are currently reading.

In this way paging allows to keep just the parts of the process that we are using in memory and the rest on the disk.

4.3.1 Basic Operation

Q. What is Page Table? Explain-Sr conversion of Virtual Address to Physical Address in Paging with example.

- The problem of external fragmentation is avoided by using paging. Paging is implemented by integrating the hardware and operating system.
- In this technique, logical memory is divided into blocks of the same size called pages.
- The physical memory is divided into fixed-sized blocks called frames (size is power of 2, between 512 bytes and 8192 bytes, also larger sizes possible in practice).
- The page size and frame size is equal. Initially all pages remain on secondary storage (backing store).
- When a process is needed to be executed, its pages are loaded into any available memory frames from the secondary storage.

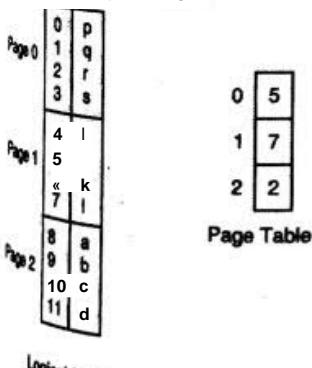


Fig. 4.3.1 : Paging model of physical memory

Memory Management

Fig. 4.3.1 : The paging model of physical and logical memory is the basic operations done in paging:

- The logical address is divided into two parts: Page number (p) and a page offset (d).
- The page number is used as an index into a page table.
- The base address of each page in physical memory is maintained by page table.
- The combination of base address with page offset gives the physical memory address that is sent to the memory unit.
- The physical location of the data in memory is therefore at offset d in page frame f.
- Because of the paging the user program sights the memory as one single contiguous space, it gives the illusion that memory contains only one program.
- But the user program is spread throughout main memory (physical memory). The logical addresses are translated into physical addresses.
- Fig. 4.3.2 shows the operation of the paging hardware.

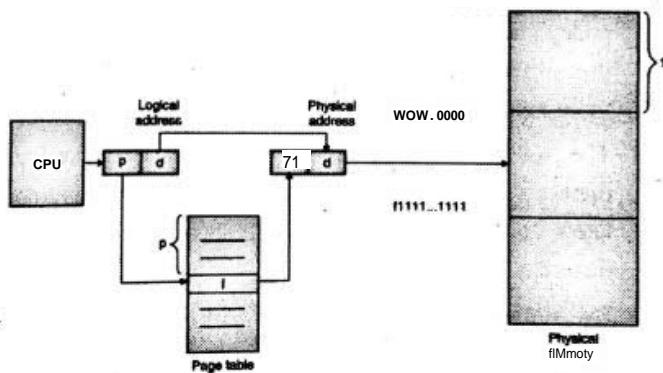


Fig. 4.3.2 : Paging hardware

- Consider the following example of Fig. 4.3.3 Let page size is of 4 K and memory available is 32 K.

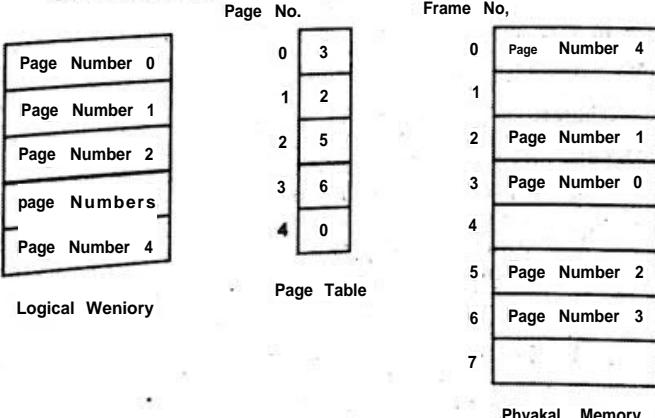


Fig. 4.3.3 : Paging example for a 32 K memory with 4K pages

TZoil frame nZbeTTpage 1 is in frame number

7 PageOisin frame number!
Logical address 4 is in page 1 and offset is 0. Page.
 $28 = ((7 \times 4) + \text{mapped to frame 7. So physical address } 28 = ((7 \times 4) + \text{offset } 15)$
 offset 0). Logical address 10 Bin page an
 Page 2 is mapped to frame 2.

So physical address $10 = ((2 \times 4) + \text{offset } 2)$.

4.3.2 Memory Protection and Sharing

- the frame number. So for the reference of the page, page table is verified to get correct frame number.
- Protection bit is associated with each page to denote whether page is read-write or read only. It is used to write read only page, a hardware operating system.
- A legal page remains in logical address space of the process and in page table this page is marked with valid (v) bit. Otherwise page is illegal and marked with invalid (i) bit.
- Paging also provides advantages of sharing of common code. The re-entrant code cannot be modified. If many users share this code, only one copy needs to be kept in main memory.
- The page table of different user points to same physical memory (frames). In this case data used by different users can be different so data pages are mapped to different frames.
- By simply having page table entries for different processes point to the page frame, the operating system can create regions of memory that are shared by two or more processes.

4.3.3 Translation Look aside Buffer

In paging, in order to access a byte first page-table entry needed to be accessed and then byte is accessed.

Therefore, two memory accesses are needed to be performed. It results in slowing down memory access by a factor of 2. This delay would be unbearable under most situations.

- The above problem can be resolved by using translation look-aside buffer which is fast searching hardware cache. The TLB is associative, high-speed memory.
- The every entry in TLB is divided in two parts one is key and other is value.**

memory is presented with an item in the associative memory at same time. If there is compared with the key in value field is returned.

match Time hardware is expensive, the search mechanism supported in this way is very fast, it holds small number of entries. Normally, TLB ranges the numbers between 1,024.

4.3.1 Effect of Page Size on Performance

→ (June 15, Nov. 15, Dec. 16)

Q. Explain effect of page size on performance.

Z?i - June 2015. Nov. 101 Deo W 5 Marks

- Important factor to decide the performance. Page size is small then any program will be divided in many numbers of pages leading to have a large page table.
- On some machine, it is required to load the page table in hardware register every time when context occurs. Hence more time will be required to load the larger size page table due to small size of the pages.
- More number of transfers to and from disk will be required due to more number of page faults.
- Hence more time will be elapsed in seek and rotational delay. If page size is large then last page will remain somewhat empty leading to internal fragmentation. Due to large size of pages, most of the unused part of program can remain in memory.
- The size of the memory divided by the page size is equal to the number of frames. If the size of page is increased then the number of frames is decreased.
- Having a fewer frames will increase the number of page faults because of the lower freedom of replacement choice.
- Large pages would also waste space by internal fragmentation.
- On the other hand, a larger page-size would draw in more memory per fault; so the number of fault may decrease if there is limited contention.
- Larger pages also reduce the number of TLB misses.

4.3.5 Hardware Support for Paging

→ (Dec-16)

Explain hardware support for paging.

MU - Dec. 2016, 10 Marks

Each operating system has its own methods for storing page tables.

Most allocate a page table for each process. A pointer to the page table is stored with the other register values (like the instruction counter) in the process control block.

When the dispatcher is told to start a process, it must **reload** the user registers and define the correct hardware page-table values from the stored user page table.

The hardware implementation of the page table can be done in several ways.

In the simplest case, the page table is implemented as a set of dedicated registers. These registers should be built with very high-speed logic to make the paging address translation efficient.

Every access to memory must go through the paging map, so efficiency is a major consideration.

The CPU dispatcher reloads these registers, just as it reloads the other registers. Instructions to load or modify the page-table registers are, of course, privileged, so that only the

The use of registers for the page table is satisfactory if the page table is reasonably small (for example, 256 entries).

Most contemporary computers, however, allow the page table to be very large (for example, 1 million entries).

For these machines, the use of fast registers to implement the page table is not feasible. Rather, the page table is kept in main memory, and a page-table base register (PTBR) points to the page table.

Changing page tables requires changing only this one register, substantially reducing context-switch time. The problem with this approach is the time required to **access** a user memory location. If we want to access location 1_{16} , we must first index into the page table, using 0_{16} value in the PTBR offset by the page number

Ex 4A1

J** Paged system TLB hit ratio is 0.9. Let the RAM access time $t = 20$ ns and TLB access time T be 100 ns. Find out

- (1) Effective memory access with TLB
- (2) Effective memory access without TLB
- (3) Reduction in effective access time.

Solution :

TLB hit ratio, RAM access time t and TLB access time T is given in problem.

(1) Effective memory access with TLB (Ea)

$$\begin{aligned} Ea &= \{(TLB \text{ hit ratio}) \times (T + t) + (1 - \text{TLB hit ratio}) \times (2T + t) \\ &= (0.9 \times (100 + 20)) + (1 - 0.9) \times ((2 \times 100) + 20) \\ &= 130 \text{ ns.} \end{aligned}$$

(2) Effective memory access without TLB (Ewt)

$$Ewt = 2T = 200 \text{ ns}$$

(3) Reduction in effective access time*

$$\begin{aligned} &= (Ewt - Ea) \times (T/Ewt) \\ &= (200 - 130) \times (100/200) = 35\% \end{aligned}$$

Syllabus Topic : Structure of the Page Table

4.4 Structure of Page Tables

→ (Dec. 14)

Q. Discuss various techniques for structuring page tables. MU - Dec. 2014, 10 Marks

Following are the most common techniques for structuring the page table.

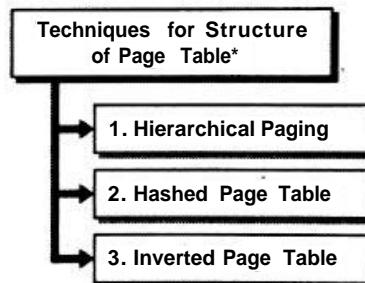


Fig. C4.2 : Techniques for structuring page tables

→ 4.4.1 Hierarchical Paging

- If the logical address space is very large such as 2^{33} or 2^{36} , the page table itself becomes extremely large.
- In this case each process may require a large physical address space for the page table alone.
- So it would not be advantageous to allocate the page table contiguously in main memory. Solution to this problem is to split the page table into smaller pieces.
- One method to achieve this is to use a two-level paging algorithm, in which the page table itself is also paged. Consider the system with 32 bit logical address space with page size 4 KB. A 32 bit logical address is divided in :
 - o 20 bit page number and

- o 12 bit page offset bove 20 bit
 - Since paging of the page tables are done, a page number is again subdivided
 - o 10 bit page number and
 - o 10 bit page offset

Thus the 32 bit logical address is divided as follows*⁵:

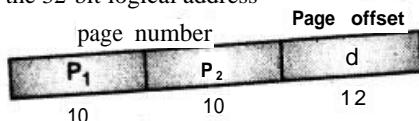


Fig. 4.4.1

where.

P, is on index into the outer page table and

P_2 is the displacement within the page of the outer page table.

- The address-translation scheme for this architecture is shown in following Hg. 4.5.2. As address translation works from the outer page table inward, this scheme is also known as a forward-mapped page table.

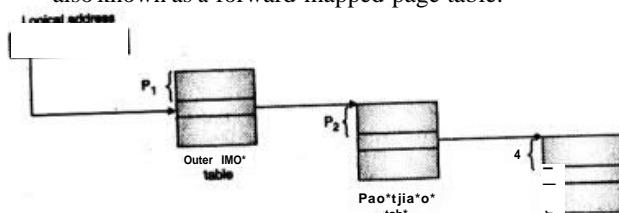


Fig. 4.4.2 : Address translation for a two-level 32-bit paging architecture

- Fig. 4.4.3 shows a two-level page-table scheme.

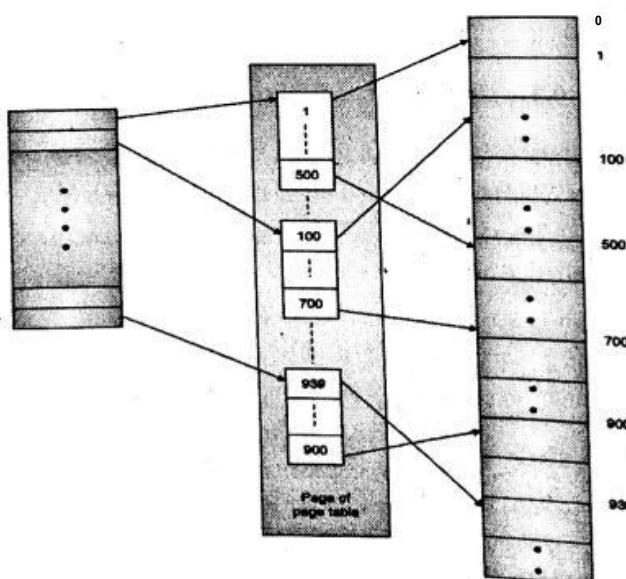


Fig. 4.43 : A two-level page-table scheme

Virtual address spaces are larger than 32 bits then page table is used in which hash value is used to find page number.

The elements which are hashed at same loc *
V a M is maintained at each entry in the ha* table
 It is recd to a v o id ** C0UiSi o n, E** elem w-u
 consists of three fields:

1. The virtual page number
 2. The value of the mapped page frame
 3. A pointer to the next element in the list

The working of the algorithm is as follows :

- Virtual address contains the virtual page number.-ft, virtual page number is hashed into the hash table.
 - The comparison of virtual page number and the field| in the first element in the linked list is performed
 - If match occurs then the corresponding page firn (field 2) is used to form the desired physical address.
 - If there is no match then succeeding entries in de linked list are searched for a matching virtual w number. This scheme is shown in Fig. 4.4.4.
 - A variant of this method that is good for 64-bit addra spaces has been proposed.
 - In this, each entry in the hash table refers to a number of pages (such as 16) rather than a single page justJilt hashed page tables.
 - Due to this, a single page-table entry can store mappings for many physical-page frames. A gnW page tables are mainly useful for address spaces,wbeff memory references are non-contiguous and scatic.J throughout the address space.

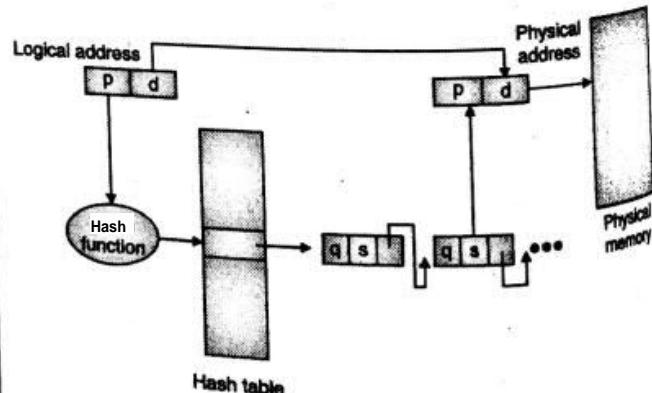


Fig. 4-4-4 : Hashed Page Table

4.4.3 Inverted Page Table

- Generally the operating system converts the logical address into a physical memory address.
- Because the table is sorted by virtual address, the operating system is capable of computing the ZT of the associated physical address entry and use that value directly.
- The disadvantages of this scheme are that each page table may contain millions of entries.
- These may use large amounts of physical memory.
- There is one entry for each real page (or frame) in memory in inverted page table. This entry contains the address of the page stored in that real memory location.
- Along with this virtual address, entry also contains the information about the process containing this page.
- Because of this only one page table is in the system, and it has only one entry for each page of physical memory. Fig. 4.4.5 shows the operation of an inverted page table.

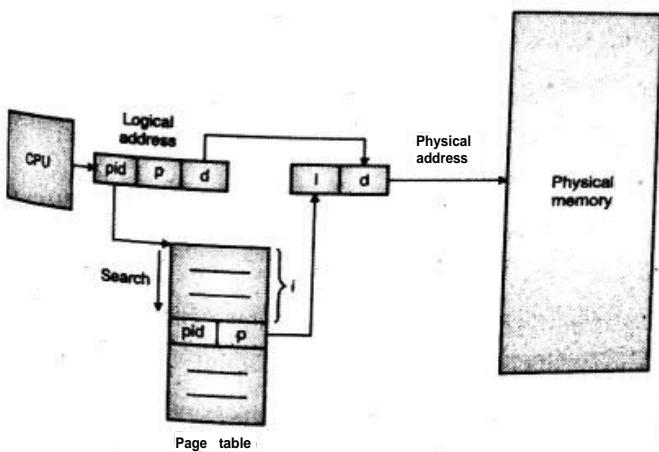


Fig. 4.4.5 : Inverted Page Table

In order to explain this approach, example of the inverted page table used in the IBM RT is demonstrated.

Every virtual address in the system is composed of a \langle process-id, page-number, offset \rangle .

As shown in above Fig. 4.5.5 inverted page-table entry contains a pair \langle process-id, page-number \rangle .

The process-id can out the work as a address-space identifier.

Memoiy Management

- At the time of the page fault, the virtual address is given to the operating system. The operating system performs a search in the page table for a match.
- If there is an illegal address produced, then there is an illegal address access has been tried.
- In this scheme, the search of the table when a page reference occurs increases the amount of time needed to search the table.
- The search would consume more time.
- To ease this problem, a hash can be used to limit the search to one or at most few page-table entries.

Syllabus Topic : Segmentation

4.5 Segmentation

Q- What is segmentation? Explain it with example.

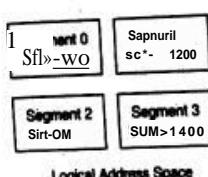
- In segmentation, a user program and its associated data can be subdivided into a number of segments. Different segments of the programs can have different length.
- Although there is a maximum segment length, the length depends on purpose of the segment in the program.
- Segments are arbitrarily-sized contiguous ranges of the address space.
- They are a tool for structuring the application. As such, the programmer or compiler may decide to arrange the address space using multiple segments.
- For example, this may be used to create some data segments that are shared with other processes, while other segments are not.
- Compiler can create the different segment for global variables, code, thread stack, library etc.
- In segmentation, a program may occupy more than one partition, and these partitions need not be contiguous.
- Internal fragmentation is eliminated in segmentation but, like dynamic partitioning, it suffers from external fragmentation. In segmentation :
 - Logical address is divided into two parts: a segment number (s) and offset (d) into that segment.

4-13



noting System (MU - Sem 4, m ==> aBSE

- Each entry of segment table contains segment base and segment limit. Segment base indicates stamping of segment in main memory and segment denotes length of the segment. Segment number is used as index to the segment table.
- X (d) always between 0 and segment limit. X will occur indicating access beyond the end of the segment. Fig. . . indicates 4 segments of the program having different sizes. Segment table contains base and limit of each segment.



	Limit	Base
0	800	500
1	1200	3300
2	WO	MOO
3	1400	1700

Segment	Base	Length
0	219	600
1	2300	14
3	1327	580
4	1952	96

Physical Memory

Fig. 45.1 : Segmentation example

Segment 3 is of 1400 bytes long and its base is 1700. Reference to byte 89 of segment 3 is calculated as $1700 + 89 = 1789$.

Example 4.5.1

On a system using simple segmentation, compute the physical address for each of following the logical address. If address generates a segment fault then indicate so.

Segment	Base	Length
0	330	124
1	676	125
2	111	99
3	498	302

(I) 0, 99 (II) 2, 7B (Hi) 1, 268 (h) 3, 222 (v) 0, 111

Solution :

Length given is limit of the segment.

- (i) 0, 99 : Segment 0 is present in segment table. Now we have to check (offset < limit) than limit then physical address = base + offset
 $99 < 124$ therefore physical address $\sim 330 + 99$

(ii) 2, 7B : Segment 2 is present in segment table. $7B < 99$ therefore physical address = $111 + 7B = 189$

(iii) 1, 268 : Segment 1 is present in segment table. $268 < 125$

X) therefore segment fault and trap to OS. Not possible

(iv) 3, 222 : Segment 3 is present in segment table. $222 < 302$ therefore physical address = $498 + 222 = 720$

(v) 0, 111 : Segment 0 is present in segment table. $111 < 124$ therefore physical address = $330 + 111 = 441$

Example 4.5.J

Consider the following segment table

Segment	Base	Length
0	219	600
1	2300	14
3	1327	580
4	1952	96

What are physical addresses of following logical addresses.

- (i) 0, 430 (ii) 1, 10 (iii) 2, 500 (iv) 3, 4000 (v) 4, 112

Solution :

(i) 0, 430 : Segment 0 is present in segment table. Now we have to check (offset < limit) or not. If offset is less than limit then physical address = base + offset

$430 < 600$ therefore physical address $\sim 219 + 430 = 640$

(ii) 1, 10 : Segment 1 is present in segment table. KK¹⁴. Therefore physical address = $2300 + 10 = 2310$

(iii) 2, 500 : Segment 2 is present in segment table. 5®¹⁴. 5⁵ (false) therefore segment fault and trap to OS. Not possible to calculate physical address

(iv) 3, 4000 : Segment 3 is present in segment table. $4000 < 580$ therefore physical address = $1327 + 4000 = 5327$

(v) 4, 112 : Segment 4 is present in segment table. $112 < 96$ therefore segment fault and trap to OS. Not possible to calculate physical address.

4.5.1 Difference between Paging and Segmentation

Q. Explain the difference between paging and segmentation. → (Dec. 0)

Mu. Dec. 2016. 5 Marks

Sr. No.	parameters	Paging	Segmentation
1	Logical address	CPU generates logical address and it is divided into two parts: a page number (p) and a page offset (d).	Logical address is divided into two parts: a segment number (s) and offset (d) into that segment.
2. Table index		The page number is used as an index into a page table.	Segment number is used as index to segment table.
3	Base address	The base address of each page in Physical memory is maintained by page table.	Segment base indicate starting physical address of the segment in main memory and segment base denotes length of the segment.
4.	Physical address	The combination of base address with page offset defines the physical memory address that is sent to the memory unit.	Segment base indicate starting physical address of the segment in main memory and segment base denotes length of the segment
5	Fragmentation	Suffer through internal fragmentation	Suffer through external fragmentation

- Memory Management**
- In segmentation:** Paging and segmentation advantages of both are:
 - me effect, \rightarrow on implementation \rightarrow 8-bit \rightarrow 0
 - In paging page also provided protection**
 - In this table:** multiple pages of the segment are divided into multiple segments. Page table is maintained.
 - Segment offset:** part of the logical address (containing offset and segment number) is further split into page offset and page number.
 - Every segment table entry includes segment base and limit of page table of that particular segment.
 - Addressing is done by page offset and page number.
 - Page Number (PN) and page offset.
 - Memory management unit uses segment number as an index into segment table in order to locate the address of page table.
 - Then PN belonging to logical address is attached to high order end of address of page table and it is used as an index to page table to locate entry in page table.
 - Frame number in page table entry is used** to calculate physical address.
 - This frame number is attached to high order end of page offset to obtain physical address.

4.6 Segmentation with Paging

Q. Explain segmentation with paging.

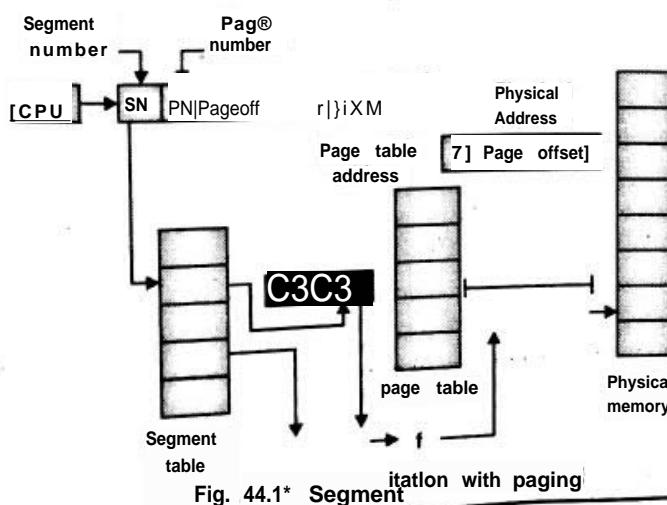


Fig. 44.1* Segmentation with paging

4.7 virtual Memory Management.

4.7.1 Virtual Memory

→ (May¹⁶)

Q. Write short note on Vttrt-1 T..J 7 ,72016. 5 Mark!

- There are many cases where entire program is not needed in main memory at a time.
- In many cases even though entire program may not all be needed at the same time.
 - Application programs always get the availability of contiguous working address space due to the idea of virtual memory.
 - Actually, this working memory can be physically fragmented and may even overflow on to disk storage.
 - This technique makes programming of large applications easier and use real physical memory more efficiently than those without virtual memory.
 - Although executing process is not entirely present in memory, it can complete its execution due to virtual memory technique.
 - Virtual memory permits execution of larger size programs although smaller size physical memory is available.
 - It means larger size programs than available physical memory can complete the execution.
 - Virtual memory concept separates the user logical memory from actual physical memory.
 - This separation offers very large virtual memory to programmers although actual physical

Syllabus Topic Tpeman agin """

4.7.2 Demand Paging

→ (Dec. 14, Dec. 16)

Q. Write short note on Demand Paging,

MJ. Dec. 2014. PeT

2016. 5 Marks

Q. What is demand paging?

- A demand paging is a paging system where pages are brought in main memory from secondary storage on demand.

is needed to execute a process, memory

W \ swapit fl * nWry process memory, memory, bto

nam J f swapping *** cn,ire memory, memory, bto

task,ead o pag ing *ll ws to SWaP " on y use pages

demand, pag ing for ejecut ion at that time. are 104

WUDI failure in accessing the page, if it j, The e dy Xped to a frame. This wUDI cause page fault which is a special type of interrupt.

- T, handle this interrupt, disk operation is initiated by S handler to read the required page in memory.
- After this mapping of page to free frame is carried out. During all this operation the process remains blocked. As soon as page becomes available, the blocked process during disk operation is now unblocked, placed on the ready queue.
- When scheduler schedules this process, it restarts execution with the instruction that caused the fault.
- The execution will continue until all the instructions are in memory.** This approach of fetching pages as they are needed for execution is called demand paging.
- Before swapping in the process in memory, the pager program makes sure about which pages will be used before the process is swapped out again.
- The complete process is not swapped in. Instead, the pager brings only those required pages into memory.
- This would restrict the swapping in of pages in memory which are not needed for execution.

0	Page 0
1	Page 1
2	Page 2
3	Page 3
4	Page 4
5	Page 5
6	Pages
7	Page?

Logical Memory

0	7	v
1		i
2	5	v
3		i
4	2	v
5		i
6		i
7		i

Page Table

0	
1	
2	Page ⁴
3	
4	
5	page5
6	
7	Page0
8	
9	

Physical Memory

Fig. 4.7.1 : Page with valid and invalid pages

- It saves swap time and the amount of physical memory needed.
- in demand paging, valid and invalid bits are used to differentiate between those pages that are in memory and those pages that are on the disk.
- Demand paging keeps more processes in memory than the sum of their memory needs.
- utilization as high as possible.

Fig. 4.7.1 shows page table with valid and invalid pages. Those pages which are in main memory are marked as valid. Pages on secondary storage are marked as invalid. Page 0, Page 2 and Page 4 are in main memory. So it is marked with valid (v) bit. Page 1 and 3 are not in main memory so marked with invalid (i) bit.

- Page 1, page 3, Page 5, page 6, and page 7 are on secondary storage.
- Continuous allocation of memory to the pages is not necessary. The reason is, my page can be mapped to any available page frame.
- The page table maintains the information regarding mapping of page to page frame. Due to this, it offers the impression of having one contiguous block of memory.

Advantages of Demand Paging

Q. List advantages of demand paging.

- Provides large size virtual memory.
- Memory is utilized more efficiently.
- Support for degree of multiprogramming is very good.

Disadvantage of Demand Paging

Q. List disadvantages of demand paging.

- Page 8th interrupt requires more number of table handling and Processor overhead with compare to simple paged memory management techniques.

4.7.3 Page Fault and Instruction Restarts

If the page is in secondary storage and not mapped to a frame in main memory, a page fault occurs. (It is a trap to the operating system)

When a page fault occurs, the operating system is responsible for:

Either killing the process in the case of an invalid memory reference, or

Memory Management

Locating the page needed into an instruction.

Following is the procedure done by operating system for handling the page fault for any process.

Internal table which is with the PCB of the process is verified for valid or invalid reference.

For invalid reference process is terminated. For valid reference, if page yet is not present in main memory

8. * in main memory from secondary

X

list .hls PaSe find emp[3/4] from frre frmc

If disk is busy operating system has to wait. If not busy search the page on disk and read it into main memory.

when reading of the page from secondary storage compares, internal table and page table of the process gets updated indicating page is now in main memory

Restart the execution of interrupted instruction from newly transferred page.

Syllabus Topic : Copy-on-Write

4.8 Copy-on-Write

- The fork() system call is used to create child process which is a duplicate of its parent. Usually, fork() system call works by creating a copy of the parent's address space for the child.
- While doing so it also duplicates the pages belonging to the parent. On the other hand, if we assume, many child processes call upon the exec() system call straight away after creation, the copying of the parent's address space may be needless.
- Copy-on-Write allows parent and child to share the same pages. If parent or child writes to shared pages then it is marked as copy-on-write and its copy is created.
- All unmodified pages can be shared by the parent and child processes.
- Pages that cannot be modified such as pages containing executable code can be shared by the parent and child. This technique used by Windows XP, Linux, and Solaris.
- As soon as it is known that a page is going to be duplicated using copy-on-write, it is significant to make

Operating System (MU - Sem 4 - IT)

a note of the location from which the free page will be allocated. Several operating systems provide free pages for such request

Syllabus Page Replacement

4.9 Page Replacement Strategies

→ (Dec. 14)

- Q. Write short note on various page replacement policies
- Q. Explain the various page replacement strategies

? After page fault main memory, it will remain always

SrSsstM--**

main memory.

For that purpose there are many page replacement strategies

- X- anheevatua ymnmn expected
- Ct num of pZ should be minimum to get the performance.
- Reference string is the string of memory references.
- The generation of the reference string is carried out artificially.
- The reference string can also be generated by listing the address of each memory reference by tracing the system.
- The important factor to determine number of page faults is the number of page frames available.
- These numbers of page faults are for a particular reference string and for page replacement algorithm for given available number of page frames.
- The number of frames available is inversely proportional to number of page faults.
- It means that if more frames are available less number of faults will occur.

4.9.1 FIFO Algorithm

→ (May 16)

- Q. Write short note on Belady's anomaly.
- Q. Compare Optimal, LRU and FIFO page replacement algorithms with illustration.
- Q. Compare FIFO and LRU page replacement algorithm.

ss == aaaa! 7r placement algorithm and wo* on the simple P³S in t (FIFO). It throws out the basis of first in t (FIFO). It throws out the order in which they were brought in, pages in the order in which they were brought with each page when it was

The 7th 8th 9th 10th 11th 12th 13th 14th 15th 16th 17th 18th 19th 20th 21th 22th 23th 24th 25th 26th 27th 28th 29th 30th 31th 32th 33th 34th 35th 36th 37th 38th 39th 40th 41th 42th 43th 44th 45th 46th 47th 48th 49th 50th 51th 52th 53th 54th 55th 56th 57th 58th 59th 60th 61th 62th 63th 64th 65th 66th 67th 68th 69th 70th 71th 72th 73th 74th 75th 76th 77th 78th 79th 80th 81th 82th 83th 84th 85th 86th 87th 88th 89th 90th 91th 92th 93th 94th 95th 96th 97th 98th 99th 100th 101th 102th 103th 104th 105th 106th 107th 108th 109th 110th 111th 112th 113th 114th 115th 116th 117th 118th 119th 120th 121th 122th 123th 124th 125th 126th 127th 128th 129th 130th 131th 132th 133th 134th 135th 136th 137th 138th 139th 140th 141th 142th 143th 144th 145th 146th 147th 148th 149th 150th 151th 152th 153th 154th 155th 156th 157th 158th 159th 160th 161th 162th 163th 164th 165th 166th 167th 168th 169th 170th 171th 172th 173th 174th 175th 176th 177th 178th 179th 180th 181th 182th 183th 184th 185th 186th 187th 188th 189th 190th 191th 192th 193th 194th 195th 196th 197th 198th 199th 200th 201th 202th 203th 204th 205th 206th 207th 208th 209th 210th 211th 212th 213th 214th 215th 216th 217th 218th 219th 220th 221th 222th 223th 224th 225th 226th 227th 228th 229th 230th 231th 232th 233th 234th 235th 236th 237th 238th 239th 240th 241th 242th 243th 244th 245th 246th 247th 248th 249th 250th 251th 252th 253th 254th 255th 256th 257th 258th 259th 260th 261th 262th 263th 264th 265th 266th 267th 268th 269th 270th 271th 272th 273th 274th 275th 276th 277th 278th 279th 280th 281th 282th 283th 284th 285th 286th 287th 288th 289th 290th 291th 292th 293th 294th 295th 296th 297th 298th 299th 300th 301th 302th 303th 304th 305th 306th 307th 308th 309th 310th 311th 312th 313th 314th 315th 316th 317th 318th 319th 320th 321th 322th 323th 324th 325th 326th 327th 328th 329th 330th 331th 332th 333th 334th 335th 336th 337th 338th 339th 340th 341th 342th 343th 344th 345th 346th 347th 348th 349th 350th 351th 352th 353th 354th 355th 356th 357th 358th 359th 360th 361th 362th 363th 364th 365th 366th 367th 368th 369th 370th 371th 372th 373th 374th 375th 376th 377th 378th 379th 380th 381th 382th 383th 384th 385th 386th 387th 388th 389th 390th 391th 392th 393th 394th 395th 396th 397th 398th 399th 400th 401th 402th 403th 404th 405th 406th 407th 408th 409th 410th 411th 412th 413th 414th 415th 416th 417th 418th 419th 420th 421th 422th 423th 424th 425th 426th 427th 428th 429th 430th 431th 432th 433th 434th 435th 436th 437th 438th 439th 440th 441th 442th 443th 444th 445th 446th 447th 448th 449th 450th 451th 452th 453th 454th 455th 456th 457th 458th 459th 460th 461th 462th 463th 464th 465th 466th 467th 468th 469th 470th 471th 472th 473th 474th 475th 476th 477th 478th 479th 480th 481th 482th 483th 484th 485th 486th 487th 488th 489th 490th 491th 492th 493th 494th 495th 496th 497th 498th 499th 500th 501th 502th 503th 504th 505th 506th 507th 508th 509th 510th 511th 512th 513th 514th 515th 516th 517th 518th 519th 520th 521th 522th 523th 524th 525th 526th 527th 528th 529th 530th 531th 532th 533th 534th 535th 536th 537th 538th 539th 540th 541th 542th 543th 544th 545th 546th 547th 548th 549th 550th 551th 552th 553th 554th 555th 556th 557th 558th 559th 560th 561th 562th 563th 564th 565th 566th 567th 568th 569th 570th 571th 572th 573th 574th 575th 576th 577th 578th 579th 580th 581th 582th 583th 584th 585th 586th 587th 588th 589th 590th 591th 592th 593th 594th 595th 596th 597th 598th 599th 600th 601th 602th 603th 604th 605th 606th 607th 608th 609th 610th 611th 612th 613th 614th 615th 616th 617th 618th 619th 620th 621th 622th 623th 624th 625th 626th 627th 628th 629th 630th 631th 632th 633th 634th 635th 636th 637th 638th 639th 640th 641th 642th 643th 644th 645th 646th 647th 648th 649th 650th 651th 652th 653th 654th 655th 656th 657th 658th 659th 660th 661th 662th 663th 664th 665th 666th 667th 668th 669th 670th 671th 672th 673th 674th 675th 676th 677th 678th 679th 680th 681th 682th 683th 684th 685th 686th 687th 688th 689th 690th 691th 692th 693th 694th 695th 696th 697th 698th 699th 700th 701th 702th 703th 704th 705th 706th 707th 708th 709th 710th 711th 712th 713th 714th 715th 716th 717th 718th 719th 720th 721th 722th 723th 724th 725th 726th 727th 728th 729th 730th 731th 732th 733th 734th 735th 736th 737th 738th 739th 740th 741th 742th 743th 744th 745th 746th 747th 748th 749th 750th 751th 752th 753th 754th 755th 756th 757th 758th 759th 760th 761th 762th 763th 764th 765th 766th 767th 768th 769th 770th 771th 772th 773th 774th 775th 776th 777th 778th 779th 780th 781th 782th 783th 784th 785th 786th 787th 788th 789th 790th 791th 792th 793th 794th 795th 796th 797th 798th 799th 800th 801th 802th 803th 804th 805th 806th 807th 808th 809th 810th 811th 812th 813th 814th 815th 816th 817th 818th 819th 820th 821th 822th 823th 824th 825th 826th 827th 828th 829th 830th 831th 832th 833th 834th 835th 836th 837th 838th 839th 840th 841th 842th 843th 844th 845th 846th 847th 848th 849th 850th 851th 852th 853th 854th 855th 856th 857th 858th 859th 860th 861th 862th 863th 864th 865th 866th 867th 868th 869th 870th 871th 872th 873th 874th 875th 876th 877th 878th 879th 880th 881th 882th 883th 884th 885th 886th 887th 888th 889th 890th 891th 892th 893th 894th 895th 896th 897th 898th 899th 900th 901th 902th 903th 904th 905th 906th 907th 908th 909th 910th 911th 912th 913th 914th 915th 916th 917th 918th 919th 920th 921th 922th 923th 924th 925th 926th 927th 928th 929th 930th 931th 932th 933th 934th 935th 936

- npare to all other page-replacement algorithms, m,
gjorithni has the lowest page-fault rate.
 - This algorithm replaces that page which will not be
ujed for the longest period of time .
 - This delays the unavoidable page fault as much as
possible, and thus decreases the total number of page
faults.
 - As algorithm replace the page on **the basis of future use
of page, it is impossible to implement.**
 - The optimal algorithm is unrealizable because it is
impossible to determine the number of instructions that
will be executed in the future before the page will be
referenced.
 - This algorithm does not suffer through Belady's
anomaly.
 - Again consider the same reference string 5, 0, 2, 3, 0, I,
3, 4, 5, 4, 2, 0, 3, 4, 3

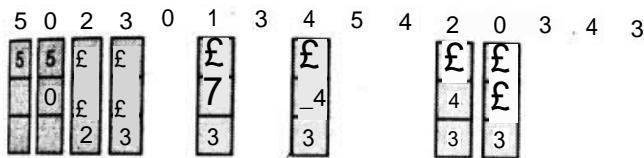


Fig.4.9.2: Optimal page replacement algorithm

- Since first three pages were initially not in main memory, references (5, 0, 2) causes page faults and brought into these 3 free frames.
 - Reference to page 3 again causes the fault. To bring Page 3 in memory as per optimal replacement policy, Page 2 is replaced because it is the page which will not be used for longest period of time.
 - Next reference is to page 0 and it is already in memory. To bring page 1 in memory again as per optimal replacement policy, page 0 is replaced.
 - Next reference is to page 0 and it is already in memory.**
 - Reference to page 4 again causes the fault. Since page 1 will no longer be referenced, it is replaced. This process continues as shown in Fig. 4.9.2 causing in total 8 page faults.

4.9.3 Least Recently Used Page

Replacement Algorithm (LRU)

- q. **Explain the replacement algorithms (LRU and HFO) with Illustration.**

Compare FIFO and LRU page replacement algorithm.

Memory Management

- The time of page's last use is associated with each Page.
 - When a page must be replaced, LRU chooses that page that was used farthest back in the past.
 - LRU is a good approximation to the optimal algorithm.
 - This algorithm looks backward in time while optimal replacement algorithm Looks forward in time.
 - **This** policy suggests that replace a page whose last usage is farthest from current time.
 - This algorithm can be implemented with some hardware support and is considered to be a good solution for page replacement.
 - This algorithm does not suffer through Belady's anomaly.

We will again consider the same reference string 5, 0, 2, 3, 0, 1, 3, 4, 5, 4, 2, 0, 3, 4, 3. Following is the result of applying LRU page replacement algorithm.

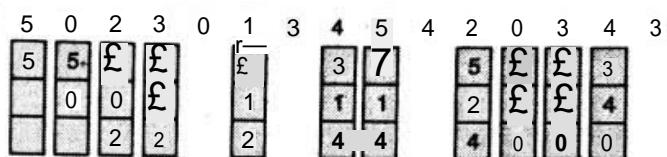


Fig. 4.9,3 : Least Recently Used page replacement algorithm

- Since first three pages were initially not in main memory, references (5, 0, **2**) causes page faults and brought into these 3 free frames. Reference to page 3 again causes the fault.
 - To bring page 3 in memory as per optimal replacement policy, page 5 is replaced because it is the page which is used least recently.
 - Next reference is to page 0 and it is already in memory.
 - To bring page 1 in memory again as per LRU replacement policy, page 2 is replaced. Next reference is to page 3 and it is already in memory. This process continues as shown in Fig. 4.9.3 causing in total 11 page faults.

< r Implementation of LRU

Using Stack

- Initially keep all the page numbers on the stack.
 - Remove the page from stack whenever it is referenced and place it on the top of the stack.
 - Any time top of stack shows latest page number that is referenced and bottom shows the page which is not used for longest period of time.

Using Counters

In a page table add time of use field with each page.

k is

- M — <> <> ■ ““ w “ 1 ““ incremented.
- Copy clock register content to time of use field
- Time of use field will show the time of last reference to the page.
- Replace the page having smallest time value.

4.9.4 LRU-Approximation Page Replacement

q, Explain LRU approximation approach.

7ZZ

4.9.4(B) Second-Chance Algorithm

The FIFO algorithm is the basic algorithm of second replacement algorithm. T: .!!!! *

- replacement, its reference bit B checked
- If the reference bit is 0 (old unused) then page is replaced.

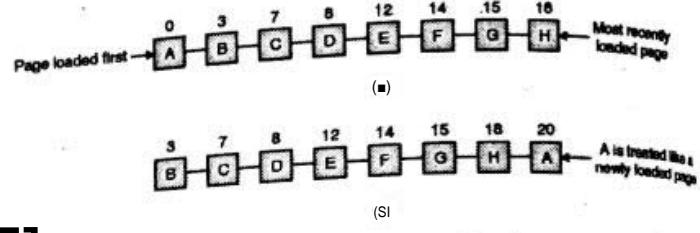


Fig. 4.9.4 : Second Chance Replacement

- Not all the system provides hardware for LRU replacement
- If hardware support not available then FIFO algorithm is used.
- The reference bit is associated with each page and stored in page table. If page is referenced for read or write then this bit is set. Initially this bit is set to 0 by OS. Once page is referenced then set to 1.
- This information is used for LRU-Approximation Page Replacement.

4.9.4(A) Additional-Reference-Bits Algorithm

- In this algorithm, 8 bit byte is associated with page and maintained in table in memory.
- After regular interval, OS is interrupted which then shifts the reference bit for each page into the high-order bit of its 8-bit byte.
- Other bits are shifted right by 1 and low order bit is discarded. This 8 bit shift register shows the history of page use in last 8 time period. For example, following values of shift register shows the use of page as shown below.
 - o 00000000 - Not used for last 8 time period
 - o 11111111 - Used at least once in each period
 - o 11000100 - This page is used more recently than one with a shift register value of 01110111
- The least recently used page is the lowest number when this shift register value is considered as unsigned integer.
- The number of bits used is hardware dependent. If single bit used then it is second chance replacement.

- If the reference bit is set to 1, a second chance is given to the page and then proceed on to select the next FIFO page.
- After a page gets a second chance, its reference bit is cleared. Now this page is considered as arrived currently and its arrival time is reset to the current time.
- Because of this, page that got a second chance will not be replaced until all other pages have been replaced (or given second chances).
- In addition, if a page is used often enough to keep its reference bit set, it will never be replaced.
- Page A was initially loaded at time 0 and given a second chance. Its arrival time is now set to current time 20.

4.9.4(C) Enhanced Second-Chance Algorithm

- This algorithm considers pair of reference bit and modify bit. We get 4 combinations with two bits as:
 - o (0, 0) - It is the best page for replacement which is neither recently used nor modified.
 - o (0, 1) - H indicates page is not good for replacement as it is not referenced but modified - need to be written out before replacement*
 - o (1, 0) - This page may be used again as it is recently used but not modified
 - o (1, 1) - This page is recently used and modified. This page may be used again and hence need to be written out to disk before it is replaced.

4.9.4(D) Clock Page Replacement

Clock \ll replacement algorithm maintains all k , frames on a circular list in the form of a clock.

The clock hand points to the oldest page.

After the occurrence of page fault, the page ~~hand~~ pointed to by the hand is checked. If its reference bit $= 0$, the page is replaced and the new page is inserted into ~~the~~ ~~lock~~ in its place.

After this insertion of new page, the hand is advanced one position. If reference bit is 1, it is cleared and the hand is advanced to the next page.

This process is repeated until a page is found with reference bit equal to 1.

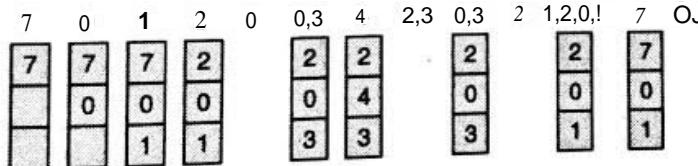
4.9.5 Counting-Based Page Replacement

4.9.5(A) Not Frequently Used or Least Frequently Used Page Replacement Algorithm (NFU or LFU)

- This policy assumes that, those pages which are referenced for more number of times will be needed again.
- A counter is used for each page and incremented with each memory reference.

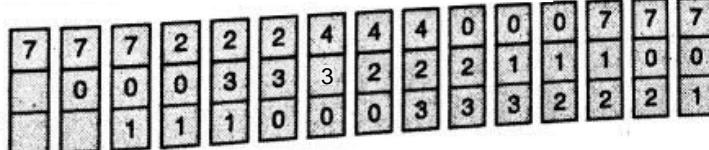
Solution:

Optimal Algorithm



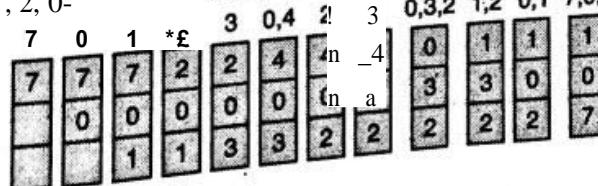
Total Page faults = 09

Optimal Algorithm: String = 7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1



Total Page faults = 15

LRU Algorithm: String = 7, 0J, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1



Total Page faults = 12

Memory Management

with less Counter value should be replaced. Suniy Se Process P has not made use of page 2 and other pages have a count of usage like 3, 4 or even 5 times.

basic ar neJd as compared to the page at 2 Se page 2 should be replaced which is having less count.

pie limitation of the algorithm is that, if initially page is used frequently and then there is no reference to this page, still it will remain in memory.

4.9.5(B) Most Frequently Used Page

Replacement Algorithm (MFU)

It is based on the assumption that, page having less count just arrived in memory and it is yet to be used.

4.10 Examples on Page Replacement

Algorithms

Example 4.10.1

For the page reference string 7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1. Calculate the Page Faults applying (i) Optimal (ii) LRU and (iii) FIFO Page Replacement Algorithms for a Memory with three frames.



Example 4.10.2

Consider the following page reference string: 1, 2, 3, 4, 2, 1, 5, 6, 2, 1, 2, 3, 7, 6, 3, 2, 1, 2, 3, 6. How many page faults occur for the following replacement algorithms, assuming four frames? All frames are initially empty, so you first

- (1) LRU replacement (2) FIFO replacement

Solution :

(1) LRU

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

Page faults for LRU = 10

(2) FIFO

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

Page faults for FIFO = 14

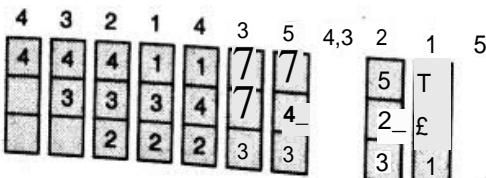
Example 4.10 J

Consider following page reference string

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

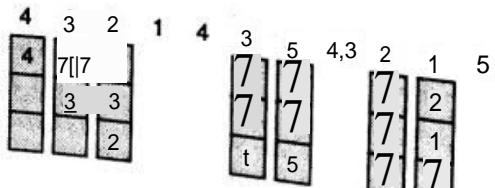
Assume frame size = 3, How many page faults would occur for FIFO, Optimal, LRU algorithm ?
Solution :

(1) FIFO page replacement algorithm



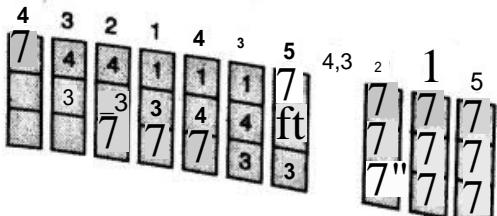
Total number of page faults = 9

(2) Optimal page replacement algorithm



Total number of page faults = 7

(3) LRU page replacement algorithm



Total number of page faults = 10

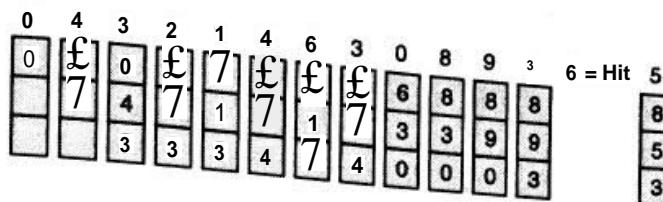
calculate Hit and Miss using (flu, optimal, FIFO page

page(ramsesi zeis3.0,4,3 2,M,6,3,0,8, 9,3,8>5 raptao *W Policies for the following sequence.

Solution :

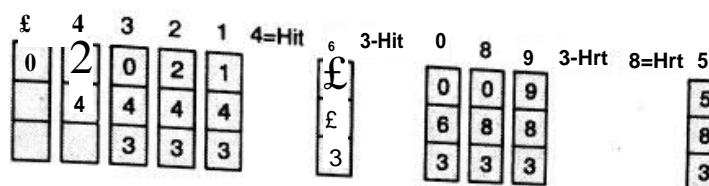
Consider frame size = 3

(1) ^{1B0}



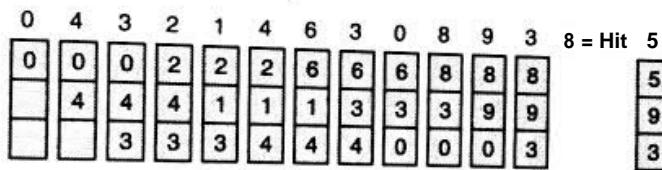
Number of hits = 01, Miss = 13

(2) Optimal



Number of hits = 04, Miss = 10

(3) FIFO



Number of hits = 01, Miss = 1

Example 4.10.5 MU - May 2016, 10 Marks

Calculate hit and miss percentage for the following string using page replacement policies FIFO, LRU and Optimal. Compare

frame size 3 and 4. Page string 2,0,3,0,4, 2,3,0,3,2,7,2,0,7,5,0,7,5,7,0

Solution:

FrameSize = 3

FIFO Page Replacement

Number of Hits - 8 Miss - 12

Frame	2	0	3	0	4	2	3	0	3	2	7	2	0	7	5	0	7	5	7	0
0	2	2	2	2	4	4	4	4	3	3	3	3	0	0	0	0	0	0	0	
1		0	0	0	0	2	2	2	2	2	7	7	7	7	5	5	5	5	5	
2			3	3	3	3	3	0	0	0	0	2	2	2	2	7	7	7	7	
PF	Y	Y	Y	-	Y	Y	-	Y	Y	-	Y	Y	Y	-	Y	-	-	-	-	

Optimal Page Replacement

Number of Hits = 13 Miss = 7

Frame	2	0	3	0	4	2	3	0	3	2	7	2	0	7	5	0	7	5	7	0
0	2	2	2	2	2	2	2	2	2	2	2	2	2	2	5	5	5	5	5	
1		0	0	0	0	4	4	4	0	0	0	0	0	0	0	0	0	0	0	
2			3	3	3	3	3	3	3	7	7	7	7	7	7	7	7	7	7	
PF	Y	Y	Y	-	Y	-	-	Y	-	-	Y	-	-	Y	-	-	-	-	-	



Number of Hits = 14 and Miss = 6

L Frame	1	2	0	3	0	4	2	3	0	2	2	2	2	0	7	5	0	7	5	7	0
0	2	2	2	2	2	4	4	4	0	0	0	2	7	7	7	7	7	7	7	7	0
1		0	0	0	0	0	3	3	3	3	3	3	0	0	0	0	0	0	0	0	7
2			3	3	3	2	2	2	2	2	2	2	2	5	5	5	5	5	5	5	5
PF	Y	Y	Y	-	Y	Y	Y	Y	-	-	Y	-	Y	-	Y	-	-	-	-	-	-

Frame Size = 4

Without Page Replacement

Number of Hits = 14 and Miss = 6

L Frame	1	2	0	3	0	4	2	3	0	3	2	7	2	2	2	5	£	7	5	7	0
0	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	0
1		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2			3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
3				4	4	4	4	4	4	4	4	4	4	4	4	4	5	5	5	5	5
PF	Y	Y	Y	-	Y	-	-	-	-	Y	Y	Y	-	Y	-	-	-	-	-	-	-

Frame Size = 4

Optimal Page Replacement

Number of Hits = 14 and Miss = 6

L Frame	2	0	3	0	4	2	3	0	3	2	7	2	0	7	5	0	7	5	7	0	
0	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
1		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2			3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
3				4	4	4	4	4	4	4	4	4	4	4	4	4	5	5	5	5	5
PF	Y	Y	Y	-	Y	-	-	-	-	Y	-	-	-	Y	-	-	-	-	-	-	-

Frame Size = 4, LRU Page Replacement

Number of Hits = 14 and Miss = 6

Frame	2	0	3	0	4	2	3	0	3	2	7	2	0	7	5	0	7	5	7	Q	
0	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
1		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2			3	3	3	3	3	3	3	3	3	3	3	3	3	3	5	5	5	5	5
3				4	4	4	4	4	4	4	4	4	4	4	4	4	5	5	5	5	5
PF	Y	Y	Y	-	Y	-	-	-	-	Y	-	-	-	Y	-	-	-	-	-	-	-

When frame size = 4 then number of hits are more all algorithms for given page string.

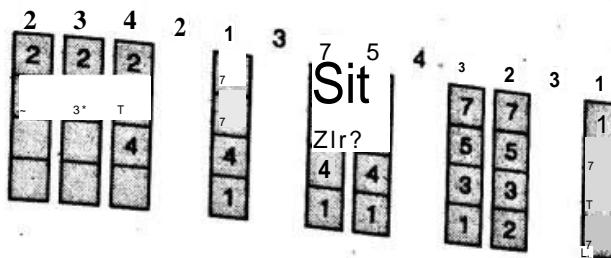
Example 4.10.6 MU - June 2015. 10 Marks

What is paging. Explain LRU, FIFO, and OPT page replacement policy for the given page sequences. Page frame size is 2, 3, 4, 2, 1, 3, 7, 5, 4, 3, 2, 3, 1

Calculate page hit and miss.

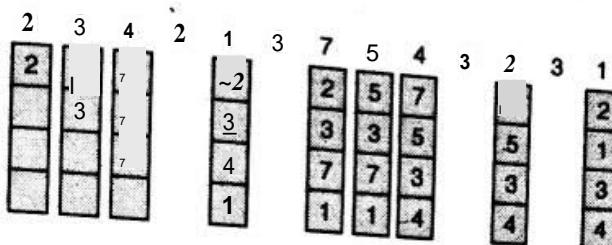
Solution :

(i) FIFO

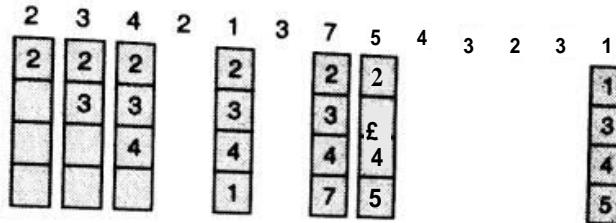


For FIFO Hit = 4 and Miss = 9

(ii) LRU

f_{LRU} Hit = 4 and Miss = 9

(iii) Optimal



fa Optimal Hit = 6 and Miss = 7

Ewnpte <<10.7 MU ~ Dec- 2016. 10 Marks

Calculate hit and miss for the following string using page replacement policies-HFO, IRC and Optimal. Compare it for the wmesize 3 & 4.

1 2 3 2 1 6 2 5 6 3 1 3 6 1 2 4 3

SoWon:

fameSi/f - 3

Wo

Frame	1	2	3	2	1	5	2	1	6	2	5	6	3	1	3	6	1	2	4	3
0	1	1	1	1	1	5	5	5	5	2	2	2	2	1	1	1	1	1	4	4
1		2	2	2	2	2	2	1	1	1	5	5	5	5	5	6	6	6	6	3
2			3	3	3	3	3	3	6	6	6	6	3	3	3	3	2	2	2	
PF	Y	Y	Y	-	-	Y	-	Y	Y	Y	Y	-	Y	Y	-	Y	-	Y	Y	Y

(iv) LRU

Frame	1	2	3	2	1	5	2	1	6	2	5	6	3	1	3	6	1	2	4	3
0	1	1	1	1	1	1	1	1	1	1	5	5	5	1	1	1	1	1	1	3
1		2	2	2	2	2	2	2	2	2	2	2	3	3	3	3	3	2	2	2
2			3	3	3	5	5	5	6	6	6	6	6	6	6	6	6	6	4	4
PF	Y	Y	Y	-	-	Y	-	-	Y	-	Y	-	Y	Y	-	-	-	Y	Y	Y

(i) upumai

Frame	1	2	3	2	1	5	2	1	6	2	5	6	3	1	3	6	1	2	4	3
0	I	1	1	1	1	1	1	1	6	6	6	6	6	6	6	6	6	2	4	4
1		2	Z	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1
2			3	3	3	5	5	5	5	5	5	5	3	3	3	3	3	3	3	3
PF	Y	Y	Y	-	-	Y	-	-	Y	-	Y	-	Y	-	-	-	Y	Y	-	-

Frame Size = 4

(i) FIFO

Frame	1	2	3	2	1	5	2	1	6	2	5	6	3	1	3	6	i	1	ft	4	3
0	1	1	1	1	1	1	1	1	6	6	6	6	6	6	6	6	I	6	6	6	6
1		2	2	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1
2			3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	2	J	2	2
3					5	5	5	5	5	5	5	5	5	5	5	5	H	Y	1	4	4
PF	Y	Y	Y	•	•	Y	-	-	Y	-	-	-	-	-	-	-	I	-	i	Y	Y

(ii) LRU

Frame	1	2	3	2	1	5	2	1	6	2	5	6	3	1	3	6	1	2	1	4	L ³
0	1	1	1	1	1	1	1	1	1	1	J	1	3	3	3	3	3	3	3	3	3
1		2	2	2	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1
2			3	3	3	3	3	3	3	6	6	6	6	1	6	6	6	6	6	6	6
3					5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
PF	Y	Y	Y	-	-	Y	-	•	Y	-	1	-	•	Y	Y	-	-	-	-	Y	Y

(iii) Optimal

Frame	1	2	3	2	1	5	2	1	6	2	5	6	3	1	3	6	1	2	4	3	
0	1	I	1	1	1	1	1	1	6	6	6	6	6	6	6	6	6	6	6	4	4
1		2	2	2	2	2	2	2	2	2	2	2	2	I	2	2	2	2	2	2	2
2			3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
3					5	5	5	5	5	5	5	5	5	5	1	1	1	1	1	1	1
PF	Y	Y	Y	-	-	Y	-	-	Y	-	-	-	-	-	-	-	-	-	-	Y	-

Frame Size = 3			Frame Size = 4			
				Hit	Miss	
	FIFO			6	14	
	LRU			9	11	
OPTIMAL		11	9	OPTIMAL		113 17
FIFO				FIFO		II 9
LRU				LRU		10 10
OPTIMAL				OPTIMAL		113 17

Example 4.10.8 MU - Dec 2015 10 Marks

Consider the following page reference string: 1, 2, 3, 4, 2, 1, 5, 6, 3, * 1, * 3, 6

occur for the following replacement algorithms, assuming three frames, LRU, FIFO, and Optimal. How many page faults would occur for LRU, FIFO and optimal page replacement.

SOW**1

Frame size = 3

LRU 15 page faults ii) FIFO, 16 page faults iii) Optimal.

The "Y" in bottom rows denote the page fault

(i) LRU

Frame	1	2	3	4	2	1	5	6	2	1	2	3	7	6	3	2	1	2	3	6
0	1	1	1	4	4	4	5	5	5	1	1	1	7	7	7	2	2	2	2	2
1		2	2	2	2	2	2	6	6	6	6	3	3	3	3	3	3	3	3	3
2			3	3	3	1	1	1	2	2	2	2	2	6	6	6	1	1	1	6
PF	Y	Y	Y	Y	Y	-	Y	Y	Y	Y	Y	-	Y	Y	Y	-	Y	Y	-	Y

(ii) H?O

Frame	1	2	3	4	2	1	5	6	2	1	2	3	7	6	3	2	1	2	3	6
0	1	1	1	4	4	4	4	6	6	6	6	3	3	3	3	2	2	2	2	6
1		2	2	2	2	1	1	1	2	2	2	2	7	7	7	7	1	1	1	1
2			3	3	3	3	5	5	5	1	1	1	1	6	6	6	6	6	3	3
PF	Y	Y	Y	Y	Y	-	Y	Y	Y	Y	Y	-	Y	Y	Y	-	Y	Y	-	Y

(iii) Optimal

Frame	1	2	3	4	2	1	5	6	2	1	2	3	7	6	3	2	1	2	3	6
0	1	1	1	1	1	1	1	1	1	1	1	3	3	3	3	3	3	3	3	3
1		2	2	2	2	2	2	2	2	2	2	2	7	7	7	7	2	2	2	2
2			3	4	4	4	5	6	6	6	6	6	6	6	6	6	1	1	1	1
PF	Y	Y	Y	Y	Y	-	Y	Y	-	-	-	Y	Y	-	-	Y	Y	-	-	Y

If frame size = 5 then

LRU gives 8 page faults ii) FIFO, 10 page faults iii) Optimal, 7 page faults

Bold page number shows the fault.

④ LRU • i, 2, 3, 4, 2, L 5, 6, 2, U 3* 3* 2, - 2 3 6

④ FIFO : j, 2, 3, 4, 2, 1, 5, 6, 2, 1, 2, 3, " * * * " 3> 6

④ Optimal. 1, 2, 3, 4, 2, 1, 5, 6, 2, 1, 2, 3, • -----

Example 4.10.9 MU - Dec 2016, 10 Marks

Calculate the page replacement for FIFO, Optimal and LRU for given reference string (assuming three frame size).

0, 5, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 5, 2, 0, 6, 6, 1

PULLON:

Frame Size 3

FIFO Page Replacement

Number of Hits = 5 and Miss = 15

Frame	6	0	5	2	0	3	0	4	2	3	0	3	2	5	2	0	5	6	0	5
0	6	6	6	2	2	2	2	4	4	4	0	0	0	0	0	0	6	6	6	
1		0	0	0	0	3	3	3	2	2	2	2	5	5	5	5	5	0	0	
2			5	5	5	5	0	0	0	3	3	3	3	3	2	2	2	2	5	
PF	Y	Y	Y	Y	Y	-	Y	Y	Y	Y	Y	-	-	Y	Y	-	Y	Y	Y	

Optimal Page Replacement

Number of Hits = 12 and Miss = 8

Frame	6	0	5	2	0	3	0	4	2	3	0	3	2	5	2	0	5	6	0	5
0	6	6	6	2	2	2	2	2	2	2	2	2	2	2	2	2	2	6	6	6
1		0	0	0	0	0	0	4	4	4	0	0	0	0	0	0	0	0	0	0
2			5	5	5	3	3	3	3	3	3	3	3	3	5	5	5	5	5	5
PF	Y	Y	Y	Y	Y	-	-	-	Y	-	-	Y	-	-	Y	-	-	Y	-	-

LRU Page Replacement

Number of Hits = 8 and Miss = 12

Frame	6	0	5	2	0	3	0	4	2	3	0	3	2	5	2	0	5	6	0	5
0	6	6	6	2	2	2	2	4	4	4	0	0	0	5	5	5	5	5	5	5
1		0	0	0	0	0	0	0	0	3	3	3	3	3	3	0	0	0	0	0
2			5	5	5	3	3	3	3	2	2	2	2	2	2	2	2	6	6	6
PF	Y	Y	Y	Y	Y	-	Y	-	Y	Y	Y	Y	-	-	Y	-	Y	-	-	-

Syllabus Topic : Allocation of Frames

4-11 Allocation of Frames

- Operating system maintains free frame list. Let us take example of single-user system with 512 KB of memory consisting of pages 4 KB in size
- This system has 128 frames each of 4 KB in size. The operating system may occupy these frames.
- 93 frames are for the user process. If we consider the pure demand paging, all pages would at the start be put on the free-frame list.
- Once user process begins execution, it would produce a sequence of page faults.
- The initial 93 page faults from the free-frame list, would all get frames

- When the free-frame list was worn out, a page replacement algorithm would be used to select one of the 93 in-memory pages to be replaced with the 94, and so on.
- If the more number of frames allocated to process, the page-fault rate decreases.
- If less number of frames allocated to each process, the page-fault rate increases, slowing process execution.
- If a page fault occurs prior to an executing instruction, the instruction must be started again. As a result, we must have sufficient frames to hold all the different pages that any single instruction can reference.

Skate ? ie> for allocation of frames

Q. Explain different strategies for allocation of frames.

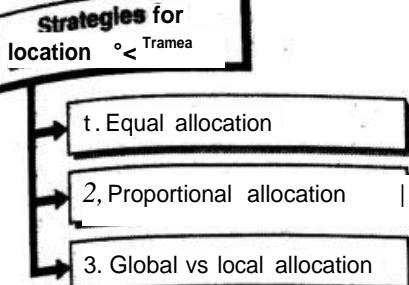


Fig. C43 : Strategies for allocation of frames

b) gqualallocHtiⁿ

Consider p number of processes and q number of frames. In this case every process will get p/q frames, to example suppose 5 processes and 93 frames present. Then $93/5$, that is each process will get 18 frames and 3 will remain at free-frame list.

* t) Proportional allocation

In this allocation, frames are allocated to the processes as per their needs. Suppose 60 free frames available. Process P1 has size 10 KB and P2 has size 128 KB. And if page size is 1 KB. Then $(10/(10+128)) * 60 = 4.34$, say nearly 4 frames allocated to P1. P2 will get $(128/(10+128)) * 60 = 55.65$, that is 56 frames.

4 X Global vs Local allocation

* Page replacement algorithm is categorized as global replacement if it can choose frame for replacement from all frames although currently it is allocated to other process also.

It means processes can take frames from other processes. Local replacement needs to replace frame by excess only from its own set of allocated frames.

The difficulty with a global replacement algorithm is that a process cannot control its own page-fault rate.

The set of pages in memory for a process depends not only on the paging actions of that process but also on the paging actions of other processes.

Consequently, the same process may perform faulty and unpredictable page faults.

Syllabus Topic : Thrashing

4.12 Thrashing

→ (June 15)

Q. What is thrashing?

MU - June 2015. 5 Marks

- A process should have some minimum number of frames to support active pages which are in memory. It helps to reduce the number of page faults.
- If these numbers of frames are not available to the process then it will quickly page fault.
- To handle this page fault, it is necessary to replace the existing page from memory.
- Since all the pages of the process are active, it will also need in future. So any replaced page will cause page fault again and again.
- Since in paging, pages are transferred between main memory and disk, this has an enormous overhead. Because of this frequent movement of pages between main memory and disk, system throughput reduces.
- This frequent paging activity causing the reduction in system throughput called as thrashing.
- Although many processes are in memory, due to thrashing CPU utilization goes low.
- When operating system monitors this CPU utilization, it introduces new process in memory to increase the degree of multiprogramming.
- Now it is needed that the existing pages should be replaced for this new process.
- If global page replacement algorithm is used, it replaces the pages of other process and allocates frames to this newly introduced process. So other processes whose pages are replaced also causes page faults.
- All these faulting processes go in wait state and waits for paging device. In this case again CPU utilization goes low.
- When CPU scheduler observes the low CPU utilization, it introduces new processes again in system to improve degree of multiprogramming. This also results in page faults of running processes.
- Again the waiting queue of paging device becomes long causing further drop in CPU utilization.
- In this way page fault rate further increases and CPU utilization decreases. There is no actual work getting done and processes spend time only in page faults.
- This thrashing can be limited by using local page placement algorithm instead of global page replacement algorithm.
- Local page replacement algorithm replaces the pages of process instead of other process. Therefore pages of the other processes will not be taken and processes will not thrash.

SlwmSetMoae!)

- Thrashing can be prevented by providing sufficient frames to the process as per **is need.**
- Locality model of process **s execution solve this problem** by keeping track on number of frames process **is the set of active pages used** currently using. I- TM
- **During execution process moves from locality to locality. Working Set Model is based on this assumption of locality. It examines most recent page references.**
- All the pages in most recent page references are the working set. Suppose working set window size defined is 4. At time t1 if the page references are **(2, 3, 2, 3, 4, 3, 3, 4) then working set is {2,3,4}** and 4 frames are required.
- At this time if demand of the process is more than working set, thrashing will occur. Operating system will monitor working set of the process and allocate needed frames.
- If sufficient frames are not available it will suspend the process. In this way working set model helps to prevent the thrashing

Syllabus Topic : Memory-Mapped Files

4.14 Memory-Mapped Files

- Q. Explain memory-mapped files.
- With compare to accessing ordinary memory, accessing files is burdensome.
 - To overcome the problem, files can be mapped to the address space of the executing process.
 - Two system calls, map and unmap, can be conceptually used for this purpose.
 - **System call causes** specified file ^{map} into a address space at a ^{the} virtual ^{system} address ^{to} ^{at}
 - Consider the file of length 128 KB which is mapped into the virtual address starting at address 1024 K.
 - Then any machine instruction that reads the byte at 1024 K gets byte 0 of the file, and the byte at 1024 K + 1 is the byte 1 of the file.
 - Similarly, a write to address 1024 K + 1000 modifies byte 1000 of the file. When PTO, terminates after finishing the execution, the modified file is left on the

disk just as though it had been changed by a combination of seek and write system calls. What actually happens is that the system's internal tables are changed to make the file become the backing store for the memory region 1024 K to 1152 K.

Thus a read from 1024 K causes a page fault bringing

to the system that supports segmentation. file mapping works better. In such a system, each file can be mapped onto its own segment so that byte n in the file is also byte n in the segment.

Suppose that this process copies file file1 to file file2. First it maps the source file, say file1 onto a segment. Then it creates an empty segment and maps it onto the destination file.

At this point the process can copy the source segment into the destination segment using an ordinary copy loop. No read or write system calls are needed. When it is all done, it can execute the unmap system call to remove the files from the address space and then exit. The output file, file2 will now exist, as though it had been created in the conservative way.

Advantage

I/O is not needed programming makes easier.

Disadvantages

- It is difficult for the system to know the exact length of the output file.
- It can happen that, if file is mapped in by one process and opened for usual reading by another and if the file is updated, that change will not be imitated in the file on disk until the page is expelled.
- The system has to take great care to make sure that processes do not see inconsistent versions of the file.
- If file is bigger than a segment, or even bigger than the entire virtual address space then alternative is to use the map system call to be able to map a portion of a file, instead of the entire file.
- Although this works, it is clearly less acceptable than mapping the entire file.

4.15 Mapped I/O

- Q. Write note on memory-mapped I/O.

- Each controller communicates with the CPU through having few registers with them.

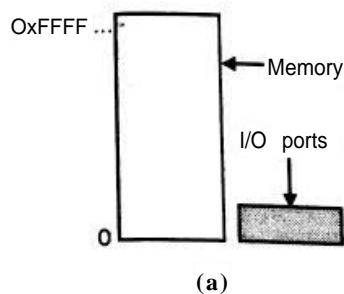
The operating system may instruct the device to send data, accept data, switch itself on or off, or otherwise carry out some action by writing these registers.

The operating system can know the device's state by reading U* ref available with device control.

This state can be: whether device is ready to accept command, and so on.

Besides the control registers, many devices are with data buffers. The operating system can read and write.

Two address



(a)

One address space



(b)

Two address spaces

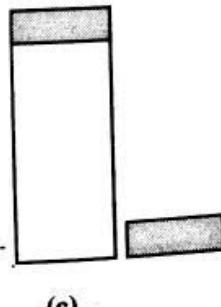


Fig. 4.15.1: (a) Separate I/O and memory space (b) Memory-mapped I/O (c) Hybrid

As an example, computers display pixels on the screen which is data. This is possible because of video RAM buffer that offers facility to store programs or the operating system for writing in it.

Memory Management

The memory management unit (MMU) interacts with the CPU and also with the device data buffers. Two approaches are available to achieve the same.

In the first method, I/O port number is allotted to each I/O port, which is an 8- or 16-bit integer.

The second approach is formed by collection of all I/O ports which is protected so that other than OS ordinary user programs cannot access it.

The second approach, introduced with the PDP-11, is to map control registers into the memory space, as shown in Fig. 4.15.1(b).

Each control register is assigned a unique memory address to which no memory is assigned.

This system is called **memory-mapped I/O**. In most systems, the assigned addresses are at or near the top of the address space.

A hybrid scheme, with memory-mapped I/O data buffers and separate I/O ports for the control registers, is shown in Fig. 4.15.1(c).

Syllabus Topic : Allocating Kernel Memory

4.16 Allocating Kernel Memory

Following are the two approaches for managing free memory that is assigned to kernel processes: the "buddy system" and slab allocation.

4.16.1 Buddy System

Q. Explain Buddy system.

The kernel produces and wipe outs small tables and buffets repetitively during the course of execution, each of which requires dynamic memory allocation. In SVR4, there is dynamic allocation of several objects by kernel as per requirement; for example these objects X proc structures, v-nodes, and file descriptor blocks. Several of these blocks are having considerably small size with compare to usual machine page size. As a result, for dynamic kernel memory allocation, the paging system seems non efficient.

In SVR4 a revised buddy system is used. Kernel memory management requires the fast process of allocation and free operations.

The disadvantage of the buddy system is the time required to fragment and combine blocks.

- Barkley and Lee at AT&T proposed a variation known as a lazy buddy system [BARK89], and this is the technique accepted for SVR4.
- The authors observed that the amount of blocks of a particular size changes slowly in time.
- of size 2 is released and is with its buddy into a block of size 2, the kernel may next request a block of size 2, splitting the larger block again
- To avoid this unnecessary combining and splitting, the lazy buddy system postpones combining until it seems likely that it is needed, and then combines as many blocks as possible.
- The lazy buddy system uses the following parameters :
 - o N_i - Current number of blocks of size 2
 - o A_i - Current number of blocks of size 2 that are allocated (occupied).
 - o G_i - These are number of blocks that are currently globally free having size 2 ; these blocks are qualified for combining; if the companion of such a block turn out to be globally free, then the two blocks will be combined into a globally free block of size 2. If standard buddy system contains free blocks (holes) then all these could be thought as globally free.
 - o U_i - These are number of blocks that are currently locally free having size 2 ; these blocks are not qualified for combining. Although the companion of such a block becomes free, the two blocks are not combined. Instead, the locally free blocks are kept awaiting future demand for a block of that size.
- The following relationship holds :

$$N_i = A_i + G_i + U_i$$

- The lazy buddy system attempts to keep a group of locally free blocks and only calls coalescing if the number of locally free blocks goes above a threshold value. If more number of locally free blocks are available, then it may happen that, at next request, the free block will not be available to fulfill the request.
- Usually, coalescing does not occur after a block, which leads to minimum accounting operational costs. Locally and globally free blocks are differentiated so loss

accordin' if locally free blocks of particular size does not exceed the number of allocated blocks of the same size ($L_i \leq A_i$)*

This is a practical rule for limiting the expand locally free blocks. ex. P_{crim} " S in \oplus ARK89i shows that this method results in remarkable savings, $D_i = A_i - L_i = N_i - 2L_i - G$

4.16.1(A) Operation of Buddy Algorithm

----- ain working operation of budtr algorithm

- The buddy algorithm manages chunk of memory as follows. Initially memory consists of a single contiguous piece, 32 pages as shown in Fig. 4.164.
- Any request for memory is first rounded up to a power of two, say 8 pages. The full memory chunk is divided in half.
- Since each of these pieces is still too large, the lower piece is divided in half again. Now suppose that a second request comes in for 4 pages. The smallest available chunk is split and half of it is claimed.

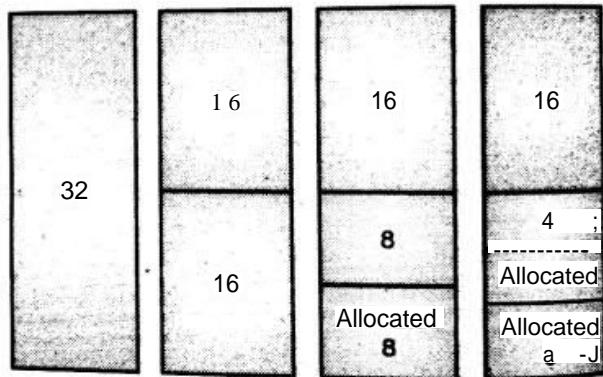


Fig. 4.16.1: Operation of the buddy algorithm

4.16.2 Slab Allocation

Q. What is slab allocation? Explain:

- Usually main memory page frames are allocated to user-space processes, dynamically allocated to ~~data~~, static kernel code, and the page cache.
- The Linux kernel memory manages physical memory page frames. Its main function is to deallocate frames for particular uses.
- memory allocation is used for user virtual memory

In the virtual memory scheme, a buddy algorithm is used to allocate memory for the kernel. The kernel can be allocated and freed in units of one or more pages.

The minimum amount of memory that can be allocated at a time is one page.

The kernel needs small short-term memory chunks in sizes of 512, 1024, 2048, 4096, and 8192 bytes. In order to hold such small chunks, Linux uses a scheme known as *slab allocation* within an allocated page.

On a Pentium/x86 machine, the page size is 4 Kbytes, and chunks within a page may be allocated of sizes 32, 64, 128, 256, 512, 1024, 2048, and 4096 bytes.

Linux also maintains a set of linked lists, one for each size of chunk. Chunks may be split and combined in a way like the buddy algorithm and moved between lists accordingly.

- The memory is allocated to kernel data structures using slab. A slab comprises one or more physically contiguous pages. A cache contains one or more slabs.
- A cache is assigned to each unique kernel data structure. For example, cache for process descriptors, file objects, and a cache for semaphores, and so on.
- The objects that are instances of the kernel data structure are occupied in the cache which represents this data structure.

This means, the cache representing process descriptors stores instances of process descriptor objects.

The relationship among slabs, caches, and objects is shown in Fig. 4.16.2.

In the following Fig. 4.16.2, two kernel objects having 3 KB size and three objects having 7 KB size are shown. These objects are stored in the respective caches for 3 KB and 7 KB objects.

The slab-allocation algorithm caches are used to store kernel objects. A number of objects are assigned to the cache after the cache is created.

716 of the associated slab decides the number of objects in the cache.

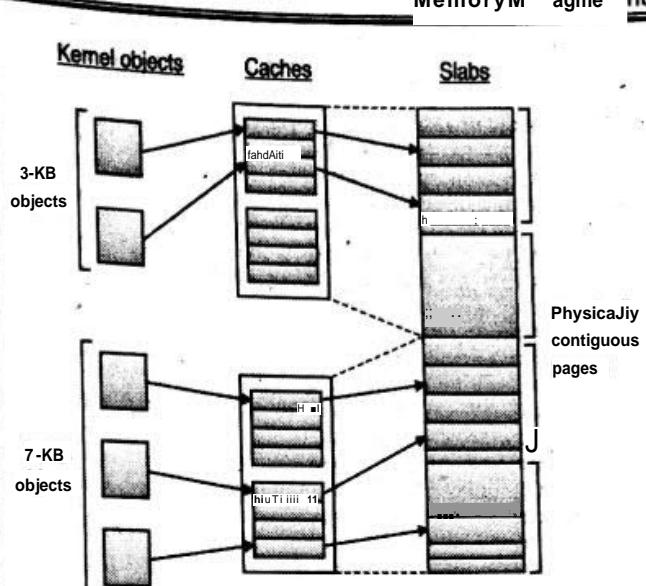


Fig. 4.16.2 : Slab allocator in Linux

- As an example, a 16-KB slab (made up of four contiguous 4 KB pages) could store eight 2-KB objects. At the start, all the objects in the cache are marked as free.
- At the time of requirement of new object for a kernel data structure, the allocator can allot any free object from the cache to fulfill the request. Then this assigned object from the cache is marked as used.

In Linux at any given point of time, a slab may be in one of three possible states :

- Full** : All objects in the slab are marked as used.
- Empty** : All objects in the slab are marked as free.
- Partial** : The slab consists of both used and free objects.

Syllabus Topic : Other Considerations

4j 7 other Considerations

4.17.1 Prepaging

Demand paging leads to many number of page faults when a process starts execution. This is due to attempt to obtain the initial locality into memory.

This can occur at other times also. For example, during restarting of swapped-out process having all its pages on the disk, and requires each page to be fetched by its own page fault.

Paging is nothing but trying to stop this high level of initial paging.



- The policy is for bringing all needed pages into memory at one time. Some operating systems, for example, Solaris prepage the page frames for small files.
- Working set model permits to keep each process a list of the pages in its working set. If it is required to suspend a process, because of an I/O wait or a lack of **free frames, the working set for that process is known.**
- For the resumption of such process due to I/O has finished or sufficient free frames have become available, it is possible to bring back into memory its **whole working set prior to its restarting.**
- Prepaging may offer a benefit in some cases. The question is simply whether the cost of using prepaging is less than the cost of servicing the corresponding page faults.

4.17.2 Page Size

- Page size is an important factor to decide the performance. If page size is small then any program will be divided into more number of pages leading to have a large page table.
- On some machines, it is required to load the page table in hardware register every time when context occurs.
- Hence more time will be required to load the larger size page table due to small size of the pages. More number of transfers to and from disk will be required due to more number of page faults.
- Hence more time will be elapsed in seek and rotational delay.
- If page size is large then last page will remain somewhat empty leading to internal fragmentation.
- Due to large size of pages, most of the unused part of program can remain in memory.

4.17.3 TLB Reach

- The hit ratio depends on the number of entries in the TLB, and the hit ratio can be increased by having more number of entries in the TLB.
- This, on the other hand, is not correct. Associative memory is both expensive and power consuming.
- **Associative memory is used to build the TLB. Hit ratio** is a metric which is similar to the TLB. Hit ratio refers to the total memory accessible from the TLB and

- is just the number of entries multiplied by the page size.
- TLB stores the working set of process. Often the process will require more time resolving references in the page table rather than the TLB.
 - If number of entries in the TLB is increased then reach also increases.
 - Other scheme to increase the TLB reach is to increase the size of the page or offer multiple page sizes.
 - It may lead to fragmentation as many applications do not need such large page size. Operating system should manage TLB for multiple page sizes instead of hardware managing it.

4.17.4 Inverted Page Tables

- The inverted page table is already discussed. The reason behind this type of page management is to lessen the amount of physical memory required for virtual-to-physical address translations.
- This saving is achieved by creating a table that has one entry per page of physical memory, indexed by the pair <process-id, page-number>.
- As inverted page table maintains the information regarding which virtual memory page is stored in each physical frame, they lessen the amount of physical memory required to store this information.
- On the other hand, the inverted page table does not include whole information regarding the logical address space of a process, and that information is necessary if a referenced page is not presently in memory.
- In order to process page faults, demand paging needs this information. An external page table should be maintained per process for availability of information.
- This each page table is similar to the conventional process page table and includes information on where each virtual page is located.
- As external page tables are referenced only once, they do not require to be available quickly. Upon a page fault, they do not require to be loaded again.
- As external page tables are themselves paged in and out of memory, they are not directly accessible by the processor. But, due to page faults, the processor has to page another page fault as it is.

table it needs to search the virtual address, in kernel this situation should be handled carefully.

4.17.5 program Structure

If user or compiler is aware of underlying paging then it is possible to improve system performance.

If data structures and programming structures carefully selected, then it can increase locality and hence decrease the page-fault rate and the number of pages in the working set.

4.18 Exam Pack (University and Review Questions)

1. Syllabus Topic : Memory Management Strategies Background

- Q. What are different requirements for memory management? Explain. (Refer section 4.1)
- Q. Write short note on Relocation. (Refer section 4.1.5)
- Q. With the help of example, clearly differentiate between logical and physical address space. (Refer section 4.1.6)

2. Syllabus Topic : Contiguous Memory Allocation

- Q. Explain memory management with Fixed and Variable Partitions. (Refer section 4.2.1)
- Q. What is internal fragmentation and external fragmentation? Explain with example. (Refer section 4.2.1)
- Q. Compare internal and external fragmentation. (Refer section 4.2.1)
- Q. What is compaction? Explain with example. (Refer section 4.2.3)
- Q. Discuss partition selection algorithm in brief. (Refer section 4.2.4) (4 Marks)

- Q. Explain different memory allocation strategies. (Refer section 4.2.4)

Example, (5 Marks) (May 2016)

3. Syllabus Topic: Paging

What is paging? (Refer section 4.3) (5 Marks) (June 2015)

Memory Management

Q. What is Page Table? Explain the conversion of Virtual Address to Physical Address in Paging with example. (Refer section 4.3) (1)

Q. Explain effect of page size on performance. (Refer section 4.3.4) (5 Marks)

(Jun 2015, Nov. 2015, Dec. 2016)

Q. Explain the hardware support for paging. (Refer section 4.3.5) (10 Marks) (Dec. 2016)

4. Syllabus Topic : Structure of the Page Table

Q. Discuss various techniques for structuring page tables. (Refer section 4.4) (10 Marks) (Dec. 2014)

5. Syllabus Topic : Segmentation

Q. What is segmentation? Explain it with example. (Refer section 4.5)

Q. Explain the difference between paging and Segmentation. (Refer section 4.5.1) (5 Marks) (Dec. 2016)

Q. Explain segmentation with paging. (Refer section 4.6)

6. Syllabus Topic : Virtual Memory Management

Q. Write short note on Virtual Memory. (Refer section 4.7.1) (5 Marks) (May 2016)

7. Syllabus Topic : Demand Paging

Q. Write short note on Demand Paging. (Refer section 4.7.2) (5 Marks) (Dec. 2014, Dec. 2016)

a. What is demand paging? (Refer section 4.7.2)

a. List advantages of demand paging. (Refer section 4.7.2)

Q. List disadvantages of demand paging. (Refer section 4.7.2)

8. Syllabus Topic : Page Replacement

Q. Write short note on various page replacement policies. (Refer section 4.9) (5 Marks) (Dec. 2014)

a. Explain the various page replacement strategies. (Refer section 4.9)

Q. Write short note on Belady's anomaly.

(Peter section 4.9) (5 Marks) (May 2016)

Q.	Compare Optimal, LRU and FIFO page replacement algorithms with illustration. (Refer section 4.9.1)	Example 4.10.9 (10 Marks)
Q.	Compare FIFO and LRU page replacement algorithm. (Refer section 4.9.1)	Syllabus Topic : Allocation of Frames (Dec 2016)
Q.	Compare Optimal, LRU and FIFO page replacement algorithms with illustration. (Refer section 4.9.2)	Syllabus Topic : Thrashing
Q.	Compare Optimal, LRU and FIFO page replacement algorithms with illustration. (Refer section 4.9.3)	Syllabus Topic : Memory-Mapped Files 1* 2015
a.	Compare FIFO and LRU page replacement algorithm. (Refer section 4.9.3)	Q. Explain different strategies for allocations. (Refer section 4.11)
Q.	Explain LRU approximation approach, (Refer section 4.9.4)	Syllabus Topic : Allocating Kernel Memory
Example 4.10.4 (10 Marks)	(Nov. 2014)	Q. Explain Buddy system. (Refer section 4.17)
Example 4.10.5 (10 Marks)	(May 2016)	Q. Explain working operation of buddy algorithm. (Refer section 4.17.1(A))
Example 4.10.6 (10 Marks)	(June 2015)	Q. Explain "hash, slab" location? Explain. (Refer section 4.17.2)
Example 4.10.7 (10 Marks)	(Dec. 2016)	
5 £!<W .8 (10 Marks)	(Dec. 2015)	

chapter

5

Module V

Storage Management

Syllabus Topics

File System : File Concept, Access Methods, Directory and Disk Structure, File-System Mounting, File Implementation, Allocation Methods, Free-Space Management, JSS, Implementation, Directory and Performance, Recovery, NFS;

Secondary Storage Structure : Overview of Mn e*

Disk Scheduling, Disk Management, RAID Structure, Disk Structure, Disk Attachment, Swap-Space Management, Storage Structure, Storage Implementation, Tertiary-Storage

I/O Systems : Overview I/O Hardware, Application I/O Interface, Kernel I/O Subsystem, Transforming I/O Requests to Hardware Operations, STREAMS, Performance.

Syllabus Topic : File System - File Concept

5.1 File System

5.1.1 File Concept

- Every application needs to access and store the information. During execution, process can keep the information in its address space. This address space to store information is sufficient for some type of applications but not for others, which requires huge size of information.

- After finishing the execution, process terminates and information also vanishes. For many applications, it is required that information should be retained forever.

- Even though computer crashes or process is killed, still information should remain retained. Most of the time, the information is sharable among many processes.

If information kept in the address space of one process, other processes will not be able to access it. It is necessary to keep information independent of the processes.

Now the necessities for long term information storage.

It must be possible to store a very large amount of information.

- o The information should not loss after termination of the process using it.
- o Several processes must be able to access the information simultaneously.
- o In order to fulfill the above requirements, it is necessary to store information on disks and other secondary storage in units called **files**.
- o A file is a named collection of related information that is recorded on secondary storage. The data cannot be written directly to secondary storage unless they are within a file. Processes can then read information from the file.
- o It also can write new information into the file if needed. After process termination, information in file should remain retained and should not vanish. A file should only vanish when its holder clearly removes it. Operating system manages the files.
- o **File system** describes how files are structured, named, accessed, used, protected and implemented.

5.1.2 File Attributes

- File is identified by its name which is string of the characters. Every file has a name and its data. All operating systems associate other information with

each file, for example, the date and time of file creation and the size of the file.

Such extra information associated by operating system is called as attributes. The lists of attributes are not same for all the systems and vary from one operating system to another.

No single existing system supports all of these attributes, but each one is present in some system.

Some of the attributes and its meaning of file

Q. Explain various file attributes in brief.

Attributes	Meaning
Protection	Access-control information determines who can do reading, writing, executing, and so on.
Size	The current size of the file
Type	Information is needed for systems that support different types of files.
Creator	ID of the person who created the file
Owner	Current owner of the file
Password	Password needed to access the file
Identifier	This unique number which identifies the file within the file system.
Hidden flag	0 for normal; 1 for do not display in listings
Read-only flag	0 for read/write; 1 for read only
ASCII/binary flag	0 for ASCH file; 1 for binary file
Lock flags	0 for unlocked; nonzero for locked
Key length	Number of bytes in the key field
Time of last access	Date and time the file was last accessed
Time of last change	Date and time the file has last changed
Maximum size	Number of bytes the file may grow to
System flag	0 for normal files; 1 for system file

Files are used to store information which can be used later on. For storage and retrieval purpose, different types of systems offer different operations. For the operation like create, write, read, reposition, operating system offer the system calls. The majority of common system calls relating to files are listed in Fig. C5.1

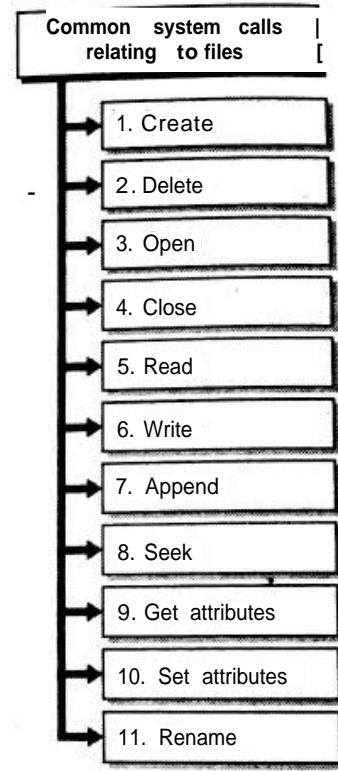


Fig. C5.1 : Common system calls relating to files

1. Create

The new file is created without data. It is needed to locate the space for this created file in file system. The intention of the call is to declare that the file is created and some of the attributes needs to be set.

2. Delete

File need to be deleted to release the disk space occupied by file when it is no longer required.

3* Open

Before making use of file, it is required to open the file first by any process. To have quick accesses in future all the attributes and disk addresses are brought into main memory.

4. Close

After use of the file completes, all accesses finished. The fetched attributes and disk address memory are no longer required, so the file should be

5.1.3 File Operations

Q. Explain various file operations in brief.

closed to release the internal table is videoed in blocks and is written to disk. Writing a file "P" ends the file's last block.

Head

The system call that is used for reading a file specifies the name of the file. Generally, the bytes are fetched from the current position. The system should maintain read pointer at the address in the file from where next read is to begin. The caller must state the amount of data desired and must also offer a buffer to put them in.

x. Write

System call - v mi. specifies file name and information to be written to the file. Writing is again at the current position. The existing file size will increase if writing starts from end of the file. In this case, the current position is at the end of file. On the contrary, if writing is done at somewhere in the middle of the file, the previously existed data will be overwritten and vanished permanently.

4.7. Append

With the help of append system call, data gets appended only to the end of the file. Systems offering less number of system calls do not generally provide append.

* i. Seek

The seek system call is used to place the pointer at a particular location in the file. It is needed for random access files to specify from where to obtain the data. After this call has completed, data can be read from, or written to, that position.

9. Get attributes

Processes frequently need to read file attributes to designated work. Get attributes system call is used for this purpose.

10. Set attributes

User can change the some attributes of file after its creation. This system call makes that possible.

11. Rename

User needs to change the name of an existing file. This call makes that possible.

5.1.4. File Naming

Process gives the name to the file when it creates it. Completion of the execution, even though process terminates, file still exists and other processes can share the same name.

Storage Management

The rules for file naming are not same for all systems, ranging from system to system, but all current systems allow strings of one to eight letters. If an operating system supports the file in logical ways, it can then operate on

Letters and special characters are also allowed to use * giving the file names. UNIX operating system differentiates between upper and lower case letters, whereas MS-DOS does not.

Windows 95, Windows 98, Windows NT and Windows 2000 support the MS-DOS file system and thus inherit its properties.

File name is divided into two parts - name and extension which are separated by period character. User and operating system can recognize the type of file from name only.

MS-DOS allows the file names containing 1 to 8 characters along with an optional extension of 1 to 3 characters. On the contrary, in UNIX user has a choice of deciding the size of extension. In UNIX, file can have two or more extensions.

Examples of file extension, Its meaning and type

Table 5.1.1 shows the some examples of file extension, its meaning, and type.

Table 5.1.1

Extension	Meaning	Type
gif	CompuServe Graphical Interchange Format image	Image
html	World Wide Web Hypertext Markup Language document	Web page
obj, o	Compiled machine language not linked	Object
exe, bin, com	Ready to execute machine language program	Executable
txt, doc	Text, document, data	Text
lib, a, dfi	Libraries of routines for programmers	Library

Extension	Meaning	Type
arc, rip, tar	Grouped files together sometime compressed for archiving or storage	Archive
tmp3, avi, mpg	Binary files containing audio or audio/video information	Multimedia
Pdfps	Its content cannot be changed	Portable document format and postscript

5.1.5 File Types

Q. Write short note on File Types.

Several operating systems support many types of files. Following are the different types of files.

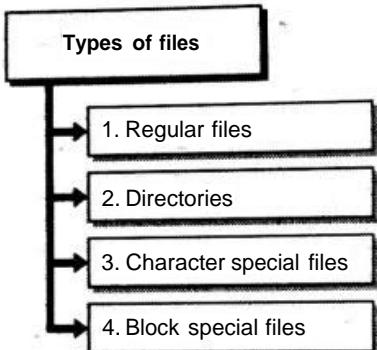


Fig. C5.2 : Different types of files

"4 1. Regular files

Regular files contain user information. The byte sequence, record sequence and tree structured explained above are the examples of regular files.

→ 2. Directories

These are system files for maintaining the structure of the file system.

→ 3. Character special files

Character special files are related to input/output and used to model serial I/O devices such as terminals, printers, and networks.

→ 4. Block special files

These are used to model disks.

- All regular files are normally either ASCII or binary files. ASCII files contain lines of text. Each line is terminated by either carriage return character or line feed character.

feed character different length. uses both types of
can be of files. Lines

The great benefit of ASCII files is that they can be displayed and printed as they are. Any text editor can edit ASCII files. Furthermore, if large numbers of programs use ASCII files for input and output, it is easy to connect the output of one program to the input of another, as in shell pipelines. Binary files are not ASCII files.

Listing them on the printer gives listing which is beyond the understanding. It gives the random garbage if it listed on the printer. Usually, they have some internal structure known to programs that use them.

UNIX binary file is an archive. It consists of, collection of library procedures compiled but w linked. **The usefulness of file types can be understood** from the TOPS-20 operating system.

The recompilation of source file will be carried out automatically after its modification. This recompilation is done if user wants to execute the object file. It guarantees that the user all the time runs an up-to-date object file. It saves the waste of time in executing the old object file.

To achieve this automatic recompilation of modified source file, the operating system must have the capability to make a distinction between source and object file. Also OS should be able to check creation time and time at which file was last modified. In order to choose the right compiler, OS must also have an ability to determine the language of the source program.

Fig- 5.1,1 demonstrate the UNIX version of execute binary file. Although this file is just a sequence of bytes, without a proper format operating system can execute this file. The file is divided into five segments. These are header, text, data, relocation and symbol table.

The very first field in header is magic number. It recognizes the type of file as an executable. It facilitates to prevent the unintentional exec* of file. Which is not in this format. After magic number, the sizes of the sections of the file are different. After all the sections of the file, the address present is the executable statement. The sizes of the address present are the flag bits are.

Program 8 the header itself. These are the test and data files and

relocated using the relocation bits. The symbol table is used for debugging.

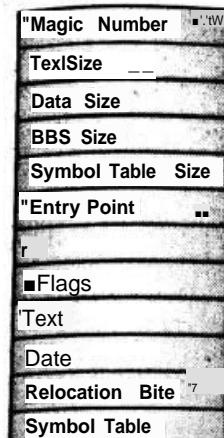


Fig. 5.1.1 : An executable file

5.1.6 File Structure

Q Explain the different techniques to structure the files.

Following are the three ways in which files can be structured.

Three ways of file structure

- 1. Unstructured sequence of bytes
- 2. Sequence of fixed-length records
- 3. Tree Structure

Fig.C53 : Three ways of file structure

5.1.7 Unstructured sequence of bytes

In this structuring, operating system treats content of the file as sequence of bytes. Any meaning must be imposed by user-level programs.

Most UNIX and Windows use this approach. This structuring provides more flexibility. Any types of content can be kept in the file and as per convenience of the user, name can be given to the file.

5.1.8 Sequence of fixed-length records

In this approach file is treated as sequence of fixed records. Each record has its internal structure. Any read operation on file returns one record and write operation will append or overwrite one record.

When punched cards were in use, record size was of 80 characters and of 132 characters proposed.

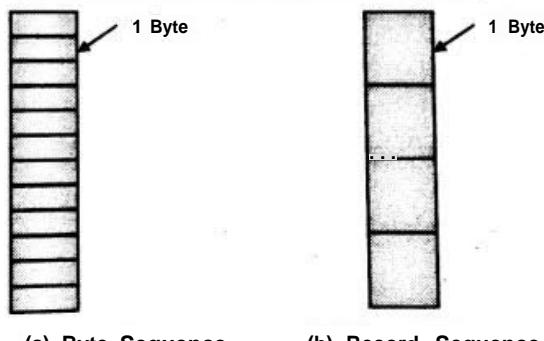
Storage Management

Printers. Many operating systems based their file systems on files consisting of 80-character records when punched cards were in use.

Programs read input from file in the unit of 80 characters and wrote in units of 132 characters keeping remaining 52 characters blank.

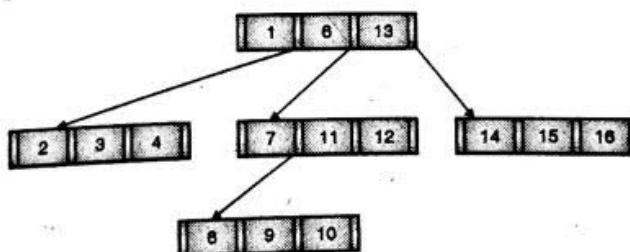
3. Tree structure

- In this organization, a file consists of the records which are organized in tree structure. All the records not necessarily have same length. Each record has a key field in a fixed position.
- The records in the tree are arranged in the sorted order of key field, to permit quick searching for a particular key.
- The main aim in this type of organization is to search record on particular key. The next record also can be searched. Operating system decides where to add new record. User is not concerned about this operation. This type of file is broadly used on the large mainframe computers and still used in some commercial data processing.
- Figs. 5.1.2(a), (b) and (c) shows above three kinds of files. Fig. 5.1.2 (c) shows the tree of records of student file arranged in sorted order of roll numbers.



(a) Byte Sequence

(b) Record Sequence



(c)Tree

Fig. 5.1J



Disadvantage of supporting multiple file structures

- If multiple file structures are supported by the operating system it would be huge in size. For example, if there are 10 file structures, the operating system would have to support 10 different file types.
- Additionally, it is required to define every file structure for one of the file type. This makes the system more complex and less efficient.

severe problems.

Opening systems like UNB, M-S-DOS enforce and support file structures which are least flexible. They treat each file as a sequence of 8-bit bytes and an operating system does not carry out the understanding of these bytes. As a result of this approach, maximum flexibility but little support is offered.

- In order to interpret an input file to the suitable structure, application should have included its own code. On the other hand, all operating systems must support minimum one structure of an executable file. Due to this support, the system will be capable of loading and running the programs.

Syllabus Topic : Access Methods

5.1.7 Access Methods

→ (June 15)

- Q. Explain different file access methods.

MU - June 2015. 10 Marks

Following are the different access methods for files.

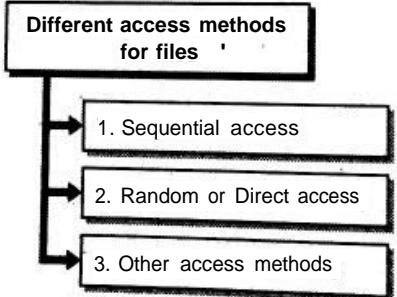


Fig. C5.4 : Different access methods for files

*4 L Sequential access

- Q. Write short note on File Access (sequential file system)

The simplest access method is sequential access. Early systems provided only this kind of file access. In this type of file access, process reads all the records in a file in order one record after other, starting at the beginning.

while accessing, skipping of any record or reading them out of order is not possible. This access method is convenient for storage medium such as magnetic tape to a certain extent than disks.

The reading of the next portion of the file is carried out in read operation and automatically file pointer is moved forward. The appending of new information to the end of the file is carried out by write operation. The pointer advances to the end of the newly write portion.

2. Random or Direct access

Write short note on File Access (Random access)

When use of disk started for storing files, it became possible to read the bytes or records of a file out of order. It is because; disks allow random access to any file block.

It also became possible to access records by its key instead of by position. Files whose bytes or records can be read in any order are called random access files. They are required by many applications such as database systems.

If railway customer calls up and wants to reserve a seat on a particular train, the reservation program must be able to access the record for that train directly instead of reading hundreds of records of other trains first.

For the random access method, the file operations must be modified to include the block number as a parameter. Thus, we have read n and write n, where n is the relative block number. Actual absolute disk address of the block is different.

The beginning of the relative block number is from address 0. Then next block number is 1 and so on although the absolute disk address of the first block is 14045 and second block is 3191.

The relative block numbers permits the operating system to take the decision about location of placing the file and facilitate the user to stop from accessing file system portions which are other than block portion.

Every read operation gives the position in the file to start reading at. In the second one, a special operation, seek, is provided to set the current position.

After a seek, the file can be read sequentially from the now-current position.

3. Other access methods

These access methods can be based on top of a random-access method. These methods generally involve the construction of an index for the file.

The index has pointers to the various blocks. The record in the file is searched by searching the index first and then uses the pointer to access the file directly to find the desired record.

Several factors are important in choice of file organization. These are

- o Minimum access time
- o Ease of update
- o Economy of storage
- o Simple maintenance
- o Reliability

Syllabus Topic : Directory Structures

5.1.8 Directory Structures

- 0. Explain different methods for defining the logical structure of a directory.
- Directories keep track of files. Directories are itself files in many systems. Systems store huge number of files on large capacity disk. In order to manage all these data, we need to organize them.
- This organization involves the use of directories.
- Following are the most common methods for defining the logical structure of a directory.

Common methods for defining the logical structure of a directory

- (A) Single-Level Directory Systems
- (B) Two-level Directory Systems
- (C) Hierarchical Directory Systems

Fig C5.5 : Common methods for defining the logical structure of a directory

→ 5.1.8(A) Single-Level Directory Systems

Structure of a single-level directory is the simplest directory structure. It contains one directory that contains all the files. This single directory is also called as a root directory.

Since on early personal computers, only one user was working, this system was more general. Fig. 5.1.3 shows single-level directory system containing five files, owned by three different users P, Q, and R. User P has two files, User Q has two files and R has one file in the directory.

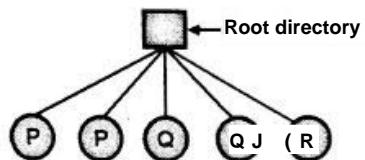


Fig- 5.13 : Single-level directory system

Advantages

- It is simple to implement.
- Locating files become faster as there is only one place to look.

Limitations

- If single user has huge number of files kept in single directory, it becomes difficult to remember the name of each file.
- If more than one user keeps the files in a single directory, then different users may give the same names to their files, violating the rule of uniqueness of names.

4 5.1.8(B) Two-level Directory Systems

- Two-level Directory System overcomes the limitations of Single-level directory. In this directory system, a private directory is given to each user. When a user refers to a particular file, only his own directory is searched.
- As different users directories are different, the same name given to the files do not interfere with each other. There is no problem in giving the same name to the files in different directories.
- Single user directory have a compulsion of having all files unique name. While creating the file for particular user, the operating system makes confirmation about whether another file of that name exists in the same user's directory or not. To know this existence of file

or not, the OS searches the directory of the user for be created. To delete a file, the has to search only directory.

zrx
Hence operating system cannot fortuitously delete another user's file that has the same name. Fig. 5.1.4 shows two level directory systems with separate users directories and its files.

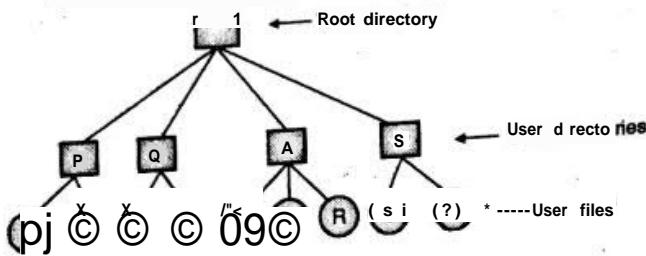


Fig. 5.1.4 : Two-level Directory Systems

- Above design multiuser computer or on a simple network of personal computers that shared a common file server over a local area network. In this system some sort of login procedure is needed.
- When a user attempts to open a file, the system knows which user it is in order to know which directory to search. Files are system programs such as loaders, assemblers, compilers, utility routines, libraries. Loader reads these files and executes when the right commands are given to the operating system.

These commands are treated by command interpreters as the name of a file in order to load and execute. In case of two level directory systems, this file name would be searched in the current user directory.

- If system files are kept in each user directory then problem can be resolved. This solution leads to wastage of memory as each user directory contains the copied system files. If special user directory containing all system files is defined the above problem can be solved.

Advantages

- Solve name collision problem as every user has separate directory.
- Independent user gets isolated from each other.

Limitations

- If the users are co-operative, that is, working on common task, then some system does not allow accessing the other user's files. In ^{some} system, if permitted, one user must have the ability to name a file

in another user's directory. To name a particular file uniquely in a two-level directory, we must give both the user name and the file name.

I. is not satisfactory oses with a k "® number of files . Even on a singled personal computer, it is not convenient.

♦ 5.1.8(C) Hierarchical Directory Systems

It is quite general and advantageous for users to group their files together in logical ^{w a Y - A} student example* might have a collection of files related different subjects of her curriculum.

First collection can be the group of files related to one subject a second collection of files related to second subject and so on. Here, some way is required to group these files together in flexible way.

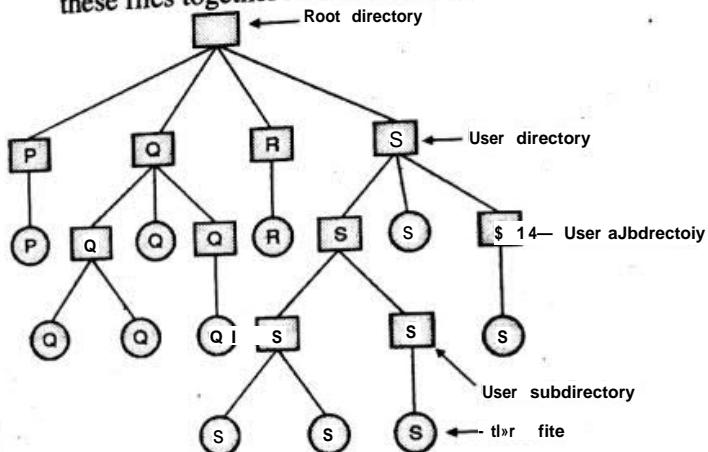


Fig. 5.1.5 : Hierarchical Directory Systems

- A tree of directories is the solution to store such group of files. A tree is the most common directory structure.
- In this approach, each user can have as many directories as are needed so that files can be grouped together in expected ways.
- The tree has a root directory, and every file ¹ the system has a unique path name. The approach is ^{shown} in Fig. 5.1.5. Root directory contains directory P, Q, R, and S which belongs to different users. Users Q and S have created the subdirectories.
- As user can create random number of subdirectories offers a commanding structuring tool for users to organize their work. This is the reason, nearly modern file systems are organized in this approach.

Advantage

With a hierarchical directory system, in addition to their files, users can access the files of other users specifying its pathname.

* 0" file System can be opened by specifying the path name. Part 8-C5.6.

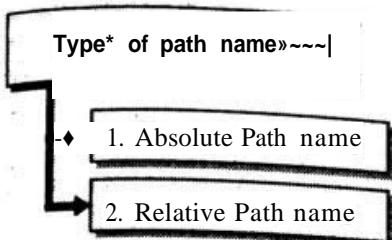


Fig. C5.6 : Types of path names

I. Absolute path name

- It consists of the path from the root directory to the file.
- The meaning of the path `/usr/myfolder/myfile` is that, the root directory contains a subdirectory `user`, which in turn contains a subdirectory `myfolder`, which contains the file `myfile`.

4 1 Relative path name

- It consists of the path from the current working directory to the file.
- A user can assign one directory as the current working directory, in which case all path names not beginning at the root directory are taken relative to the working directory.
- If the current working directory is `/usr/myfolder`, then the file whose absolute path `/usr/myfolder/myfile` can be referenced simply as `myfile`.
- In UNIX operating system the elements of the path are separated by `/`. In Windows the separator is `\`. The same path `/usr/myfolder/myfile` in UNIX, is specified as `\usr\myfold\myfile` in Windows.
- Path separator then path is absolute path. If during program always needs a particular file, it should use absolute path to access that file from any working directory.

5.1.10 directory Operations

Q/ 2' Vw different directory operations To demonstrate the directories, different system calls with respect to more dissimilarity from system to system with Wm calls for files.

idea of adding of system calls for Consider the following examples from IT.

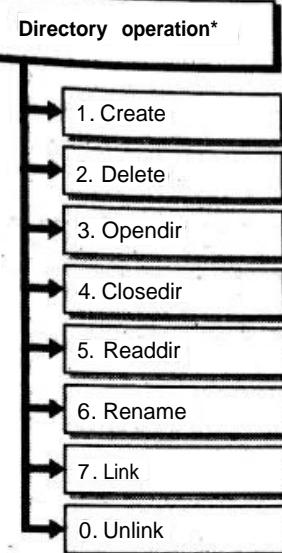


Fig. C5.7 : Directory operations

1. Create

Empty directory is created excluding dot and dotdot (.) and (..), which are placed there automatically, by the operating system,

2. Delete

If only single dot and double dot (.) and (..) is present in directory then it is treated as empty, single dot and double dot within directory cannot usually be deleted.

3. Opendir

Directory can be read to list all the files from it. Directory should be opened before reading just opening and reading the file,

4. Closedir

After reading completes, a directory should be closed to free up inner table space.

5. Readdir

It returns the next entry in an open directory, `readdir` always returns one entry in a standard format irrespective of possible directory structures is being used.

6. Rename

Directories can be renamed just like files. `Rename` renames the directory.

7. Link

Due to linking, a file appears in more than one directory. `Link` system call specifies an existing file and



its path name, and creates a link from the existing file to the name specified by the path. In this fashion, the same file may come into view in several directories.

8. Unlink

Unlink removes the directory entry. If the file unlinked only exist in one directory then it is deleted from file system.

Syllabus Topic : File System Mounting

5.1.11 File System Mounting

As discussed earlier, a ms* can be divided in multiple partitions and each partition can have its own file system or there can be multiple disks connected to system. So there can be multiple file systems. It is possible to copy a file from one file system another. Consider the disk-1 and disk-2 with file systems as shown in Fig. 5.1.6 circles shows files and rectangles shows directories.

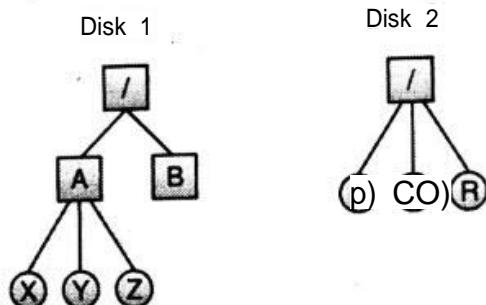


Fig. 5.1.6 : Two file systems on two disks

Suppose we have to copy file R from disk-2 to a directory A under disk-1. After mounting directory structure becomes as shown in Fig. 5.1.7.

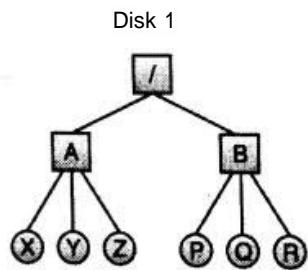


Fig. 5.1.7 : Directory structure after mounting

Now the command `cp /B/R/A/R` copies the resultant file R. In UNIX a facility is given to mount or unmount a file at prespecified time. The mount and unmount system calls are used for this purpose. A system administrator does the mounting at the time of booting of UNIX. A large number of file system are supported by UNIX for mounting.

Implementation of mount/unmount

File systems can be on same or different devices. With the mount these file system can be viewed as single file system. Files can be moved between different director of different file systems thereby appearing as a single system. The mounting can be realized in UNIX through mount table. Along with other information, the mount table contains following information.

- I-node of root node of the file system to be mounted.
- I-node of the directory in which mounting is to be carried out. This directory is called as mount point.

The mount table allows the kernel to traverse from one file system to other. The user remains unaware about it under the feeling of single file system.

Implementation of link/unlink

If we need to access file P12 from P1 and also from R1

the Link needs to be created from R1 to P12. To this shared file now two absolute pathnames exist to this shared file as shown in Fig. 5.1.8. These are /P/P1/P12 and /R/R1/P12. In directory R1 the same file can be called by different name say R1. In this case, pathnames /P/P1/P12 and /R/R1/R1 would be the same.

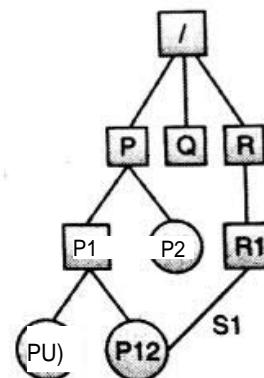


Fig. 5.1.8 : Linking

Consider the file to be linked and therefore shared is S1 in directory R1. The same file is P12 in P1 as shown in Fig. 5.1.8. The linking is done by kernel with link system call as follows.

- Except super user, other users will be not permitted to execute link system call. Error messages will be outputted if user is not super user.
- The path name /P/P1/P12 is parsed to obtain the i-node of the file to be linked. Store it for future use.
- i-node record is accessed and number of links is incremented.

in order to have disk consistency, the i-node ~~s write~~ back to disk.

The path name /R/R1 is parsed to obtain the i-node ~~no fe~~ of the directory. The content of the directory is then *read free* entry in directory is located and it makes sure that file with name S1 not present in directory. If ~~pre ent~~ message is outputted and exit operation is performed.

The entry of the Jinked file is done in /R/R1 with S1 as file name and i-node number stored in step 2.

all directories are then written back.

5.1.2 Working of Fifes

While working on a common task, many users need to share their files. If these shared file appears simultaneously in different directories belonging to different users, then it will be convenient to all users.

Fig. 5.1.9 shows that one file of user S appears in the **directory of user Q. Because of the sharing, file system is now Directed Acyclic Graph (DAG) instead of tree.**

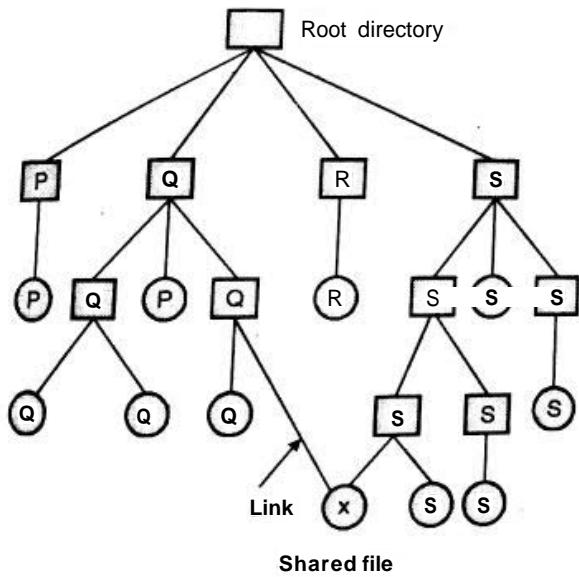


Fig.5.1.9 : File system combining a shared file

Problem in sharing the files

Suppose directories actually do contain disk addresses then, in above example, copy disk addresses of S's directory should be made available in Q's directory when the file is linked.

The user who actually perform append operation, the appended block appears only in his directory* In ~~above~~ example, If either Q or S perform append operation to

the file, then the new blocks will be listed only in the directory of the user, who is actually doing the append. The changes will not be visible to the other user, thus overwhelming the intention of sharing. Following are two methods to handle this problem.

- Instead of listing the disk blocks in directories, it is listed in data structure related to the file, just like i-node in UNIX. The directories would then point just to this data structure.
- When Q link to S's file, system create a new file of type LINK and enters this file in Q's directory. The new file contains the path name of the file to which linked. In our example, it is S. When Q read from linked file, the operating system note that it is type LINK file and search the name of the file and read it. This is called symbolic link approach.

The disadvantage of first method is that, when Q links to the shared file of user S, i-node records owner of file as S and link count becomes 2 which was initially 1. Now system knows 2 directory entries pointing to the file. If S removes the file and clears the i-node, then Q's directory entry will point to invalid i-node. If later on i-node is assigned to other file then Q will point to the wrong file.

- Due to link count in i-node, operating system will come to know that file is in use but unable to find all the directory entries of that file. The solution to this problem is to remove S's directory entry, but do not clear the i-node and keep it with count set to 1.
- Now Q is the only user having a directory entry for a file owned by S. Here system will assume that file belongs to S until Q decides to remove it. When file is deleted, count goes to 0 and no longer belongs to any user.
- The disadvantage of symbolic link is that, extra operating cost required. The LINK file having the path must be read, and then parsing of path is done and followed, component by component, until the i-node is searched.
- All of this activity may require a considerable number of additional disk accesses. Every symbolic link requires additional i-node, as is an additional disk block to store the path.
- Although path name is short, the system could store it in the i-node itself for optimization purpose.

Sullabun Topic: Sharing

5.1.13 File Sharing

Q. Write note on 'access rights'

It is necessary to share a file among multiple of users in multiuser system. It needs to deal with two issues: that are access rights and the management of simultaneous access.

5.1.13(A) Multiple Users

Q. Explain issues related to file sharing.

Q. Explain different access rights to the file.

The file system should offer a flexible facility for permitting widespread file sharing among multiple users. It offers a many system to

is also a responsibility to the way of file accessing. In alternatives in order to control access to files. In general, it is important to grant partly access a particular

file. A user with less rights has been used. The following are the some of the examples of access rights that can be granted to a particular user for a particular file.

Examples of access rights

- 1. None access right
- 2. Knowledge
- 3. Execute
- 4. Read
- 5. Append
- 6. Update *
- 7. Changing protection
- 8. Delete

Fig. C5.8 : Examples of access rights

→ 1. None access right

The user may not still know about presence of the file.
To put into effect this access limit, the user is restricted from reading the user directory that contains this file.

→ X Knowledge

With knowledge access right, the user can find out **existence of file and owner of it. The user can appeal the owner for added access rights.**

→ 3. Execute

Program is possible to user A loading and execution of program is possible to user but copying it is not allowed.

→ 4. Read

A read access right allows the user to read the file for C.y work, together with copying and execute

→ 5. Append

The user can only add data to the file at the end.

→ 6. Update

The update access right allows the user to modify, delete, and add the data to the file.

→ 7. Changing Protection

Changing protection the access rights assigned to the file. Normally this right is associated only with the owner of the file. The owner is allowed to extend this right to other users in some systems

→ 8. Delete

The user is allowed to delete the file from the file system as per need.

These rights assigned to user form a hierarchy, with each right involving those that precede it. Thus, if any user is assigned the read right for a particular file, then that user also gets following rights : knowledge, execute. The person able to create file is by default the owner of that file.

The owner is granted all of the access rights and may put rights to others.

The different categories of users to which access rights are given are :

- Specific user : Individual users having its own user ID.
- User groups : A group of users who belong to a particular user group. The system keeps track of membership of user groups.

- AU: Authenticated users of this system/

< Simultaneous Access

It is necessary to protect the share simultaneous updating from more than one user. If more than one user granted access to append a file, the operating system or file manager must implement some way to restrict it.

A brute-force scheme permits a user to complete file when it is to be updated, it is also possible to lock individual records during update.

Explain in detail.

a

Sharing of data in the form of file communication over the network is implemented is transfer using ftp.

The second method is distributed file system (DFS), it stores directories are abject to be seen from local machine. In World Wide Web (WWW) is required to get access to the remote essence a wrapper for ftp is used to supports both anonymous and authenticated access.

In anonymous access, users do not have remote machine, they can transfer file? makes use of anonymous file almost exclusively.

A tighter integration is required in DFS between the machine accessing the remote files the machine providing the files.

* The Client-Server Model

Remote file system permits mounting of one or more file system from one or more remote computers. Machine storing the file is server and machine from which file is accessed is client.

A server can offer the service to multiple clients, and a client can request to multiple servers. This depends on how client-server facility is implemented.

The server generally states the available files on a volume or directory level.

Client is specified by network name or IP address. As these can be spoofed, unauthorized user can access the server.

Secured authentication is required using encrypted keys. During exchange of these keys, unsecured access is also possible.

UNIX and its network file system (NFS), the user IDs on the client and server must match for authentication. In mismatch case, the server will be able to decide access rights to files.

Chuted Information System

Distributed naming service is also used to support distributed access to information Putation.

remote file system.

For internet (WWW) domain storage management is used. Wide use of machines in network. Undoubtedly, many types of methods are used as network known as network. Microsystem known as network. Industries adopted its names printer information, and others. But NIS is more secure than NIS.

prefers the use of Light weight directory naming method.

Active directory is based on light weight directory with the OS and is used for user authentication in Edition to system-wide retrieval of information, for example availability of printers.

Organizations can use single distributed LDAP to store all user and resource information for all their computers. This allows secure single sign on for all the users.

Failure Modes

As there can be failures in local file system due to various reasons, there are more failure modes for remote file system. This is due to the interference in communication over network and complexity of network systems.

There can be interruption in network due to hardware failure, poor hardware configuration, or networking implementation issues.

A single failure can interrupt the flow of DFS commands.

Due to crash or failure in server remote file system will be unable to reach. In this case, system can either terminate all operations to the crashed server or delay the operations until the server is again reachable.

Implementation of the failure semantics is a part of the remote file system protocol.

Loss of data can be resulted due to termination of operation. Hence delay is preferred.

Recovery from failure requires maintaining the state information at both client and server. If server must know that it has remotely mounted exported file

systems and opened files. NFS use stress DFS
But this stateless DFS is not secure.

5.1.13(C) Consistency Semantics

File systems supporting file sharing are evaluated
important criteria which is consistency semantics. If one
user makes modification in file, then other user
able to see it.

UNIX Semantics

IX

Following consistency semantics are used by the UN
file system.

ew

- Other users who have file open, can immediately view writes to this open file by a user.
- One mode of sharing permits users sharing of the **pointer of current location into the file**. Hence, If one user advances the pointer then it affects all sharing users. Here, a file has a single image that interleaves all accesses, in spite of their origin,

Session Semantics

Following consistency semantics are used the Andrew
file system (AFS).

- Other users who have file open, can not immediately view writes to this open file by a user.
- Once a file is closed, the modifications done in it are visible only in sessions starting later on. Already open instances of the file do not reflect these modifications.

Immutable-Shared-Files Semantics

For immutable-Shared-Files, if creator of file declares it as shared then changes cannot be made in it. The name of an immutable file may not be reused and its contents may not be changed. Therefore, the name of an immutable file indicates that its contents are fixed.

Syllabus Topic : Protection

5.1.14 Protection

Q. Explain system protection in detail.

- Protection mechanisms refer to the particular operating system mechanisms which are used in managing information, files and resources in the computer.
- Policy means whose data should be protected from misuse. **mechanism is how system put into effect these policies.**

In some systems, a program called reference monitor is used to impose a protection. Any attempt to access the resource is verified by reference monitor whether it is legal or not. The reference makes a decision on the basis of policy table.

Modern protection mechanism has developed to improve the reliability of any complex system if many resources are involved. There are reasons to offer the protection.

It is needed to provide the protection to avoid deliberate violation of an access restriction by a user. It is necessary to make sure that each program component running in the system uses system resources as per defined policies to ensure the reliability of the system.

Policies to make use of the resources of the system are put into effect by mechanism which is offered by protection. Some of the policies are included in the design of the system. Other policies are decided by the management of a system. Many policies are defined by the individual user of the system to protect their own files and programs.

System should be flexible enough to offer different and put into effect the different types of policies. Different types of applications have different types of resource use.

The policies designed should allow offering this resource use need. The resource use of the applications can change over the period of time. Hence, instead of relying totally on operating system, application programmer should use protection mechanism to protect the resources created against misuse.

5.1.14(A) Type of Accesses

Protection mechanisms give controlled access by restricting the types of file access that can be done. Access is permitted or denied because of a number of factors. An example of this is the request type. Following are the type of operations that may be controlled :

Read : Indicates reading from the file.

Write : Indicates writing or rewriting the file.

Execute : Indicates file must be loaded into memory and then carry out execution of it.

Append : Write new content at the end of the file.

Delete : Deletion of file and reuse the freed space.

List : List the name and attributes of the file.

stations, example rena J
or editing the file, may also be controlled. In many systems, copying, and operations are implemented as system calls. For example, copying a file from one location to another is implemented by a process. Hence, a user having read access can also cause it to be copied, printed, and so on.

5.14(B) Protection Domains

The system contains many hardware objects such as memory segments, disk drives, printers, and many software objects such as files, semaphores, and processes. Each object can be referenced by every object in the system. For example, WAFT and SIG semaphores and READ and WRITE on files.

A system should enforce a mechanism, to restrict the processes from accessing the needed objects for which they are not authorized. The mechanism should also ensure to restrict processes to a subset of the legal operations when that is needed.

For example, process P has a permission to read file F but not of performing write operation on it. A set of object and rights pair is called domain. Each pair denotes an object and some subset of the operations that can be performed on it. One domain corresponds to one user and specifies the permissions to user for certain activities. Consider the Fig. 5.1.10 three domains.

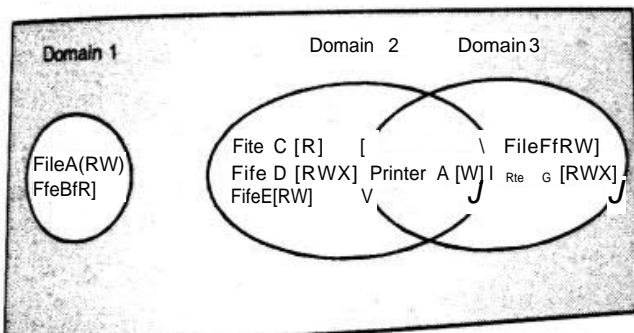


Fig. 5.140 : Domains

It is possible for the same object to be in multiple domains, [Read, Write, Execute] rights, each on each object. At a particular time of execution, a process executes in some protection domain. It can be in that domain there is some set of rights for access, and for each object it has some shown in square brackets.

During execution, processes can go from one domain to another domain. The rules are very much depends on and asterisk.

Storage Management

In Unix every object is represented by a file ID (FID). The file ID contains the file name, file type, file size, file owner, file group, file permissions, and file creation time. The file ID is used to identify the file and to access it. The file ID is also used to protect the file from unauthorized access.

5.1.14(C) Access Control

- The kernel maintains a list of objects from the user's point of view. The system keeps track of which object belongs to which domain. The access matrix for above three domains is shown in Fig. 5.1.11. The matrix tells if an access to a given object in a specified manner from a specified domain is permissible.
- For this purpose, the system makes use of this matrix and the current domain number.

		Object							
		File A	File B	File C	File D	File E	File F	File G	Printer A
Domain	1	Read	Read						
	2			Read	Read write Execute	Read Write			Write
3						Read Write	Read write Execute	Write	

Fig. 5.1.11 : Access Matrix for Fig. 5.1.10

- It is possible to switch domain in the matrix model by considering that a domain is itself an object, with the action enters.
- processes in domain 1 are allowed to go in domain 2, but once it is in domain 2, they cannot go back to domain 1.
- This situation models executing a SETUID program in UNIX. In this example, other domain switches are not permitted.

rpomain/		File A [File &) File C I		File D	File E	File F	File G	Printer A	Domain 1	Domain 2	Domain 3
1	Read Write	Read									Enter
2			Read	Read write	Read write				Writ*		
3						Read Write	Read Wn te Execute	Write			

Fig. 5.1.12: Matrix ^{nMd} *| wttli<lo"“ -ns switching

v Disadvantages of access matrix

- In fact storing the matrix is not preferred practically due to its large size and it is also sparse.
- **Most domains have no access at all to most objects, so** storing a very large, mostly empty, matrix is a waste of disk space.

The two practical methods to store the matrix are, first storing the matrix by rows or by columns, and then storing only nonempty elements. In the first method, an ordered list consists of all the domains that may access the object, and in what way. This list is associated with each object and called as Access Control List (ACL). In the Fig. 5.1.13 shown below, three processes, each belongs to different domains X, Y, and Z and three files F1, F2, and F3 are shown. Each domain corresponds to exactly one user, X, Y, and Z.

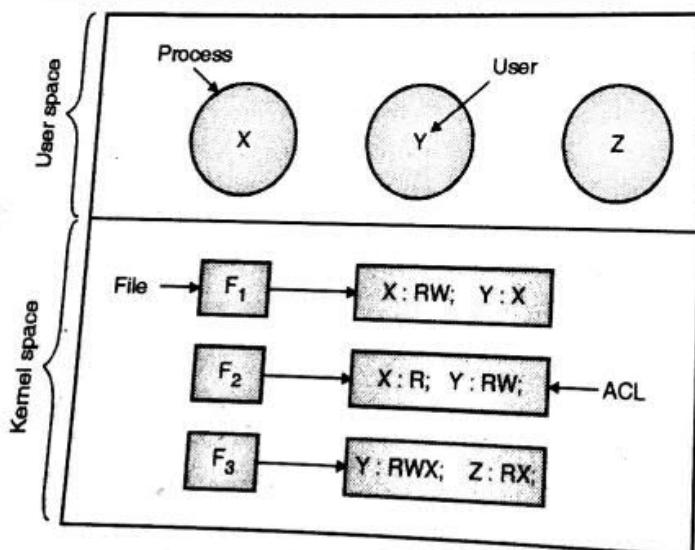


Fig. 5.1.13 : Access control list

- File entries in ACL are separated by semicolon. Each file has an ACL coupled with it. File F1 has two

entries in to ACL. The first entry indicates that any process begged by user X may read and write the file. The second entry indicates that any process Xngetouserrmayreadthefile.

Other than these accesses by these users are prohibited. Apart from this, all accesses by other users are also prohibited. Again, process does not approve the rights, these are approved by user. Any numbers of processes belonging to user X can read and write file F1. It is the owner who plays the role for protection not the process ID.

ACL contains three entries for file F2 : X, Y, and Z can alt read the file, and also Y has a right to write it. Apan from these accesses other accesses are not permitted. File F3 seems to be an executable program, since Y and Z can both read and execute it. Y can also write it.

Above example demonstrates the fundamental form of protection with ACL's. More sophisticated systems are frequently used in practice. Other than read, write and execute, many other rights are given in practice. Apart from, having ACL entries for single user, many system supports group of users. Groups have names and can be included in ACLs.

Sy *biiS-Topic File System Implementation

e System Implementation

Secondary storage such as disk is designed to hold data permanently. The file system resides in the file system. A secondary storage. Implementation of file system ideals with the following.

In what way files are stored in directories and management of the files & storage is managed.

that K is to ensure that everything works.

unfailingly work efficiently

pushes off the mass of secondary storage.

file system is maintained. The data is stored in main memory > performed in unit of bytes between disk.

*Lii> one or more sectors. These sectors are stored in disk which may be 4,096 bytes depending on disk.

sector is 512 bytes. Because of drives used. Usual

J is suitable medium for storing various files.

A block can be read from disk and modification can be written back in the same place.

Any given block of information can be accessed by disk which makes it simple to access sequentially or directly, and switching.

sorter requires only moving the read/write head to

waiting for the disk to rotate and

5.2 J File System Structure

[Q] Explain file system structure?

- The file system contains several various levels. The layered design of file system is shown in Fig. 5.2.1.
- In this design, each level in the design uses the features of lower levels to create new features for use by higher levels.
- The lowest (bottom) layer consists of device drivers and interrupt handler to exchange information between main memory and disk.
- The file-organization module keeps information about files, their logical blocks and physical blocks.
- The file organization module can translate logical block addresses to physical block addresses for the basic file system to transfer. This is achieved by identifying the type of file allocation used and the location of the file.
- Logical block numbers for each file's block is numbered from 0 or 1 to N.
- The data is stored in physical blocks. These physical blocks do not match the logical numbers. Hence, a translation is required to find each block.

Free-space manager is also a part of file-organization module. It tracks unallocated blocks and keeps these blocks to the file-organization modulo when

Storage Management

physical management logical flip metadata information is carried

Metadata carries all of the file-system structure. It

The management of files maintains actual data or contents of the files.

by logical file system in order to off the file

paired, given the symbolic name. It maintains file

Return via fit. name. It maintains file

FCS is X

information about UNIX

permissions, etc.

The

security is also protection

The

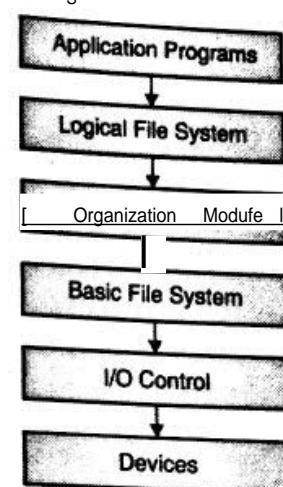


Fig 5.2*1: Layered File system

Advantage of layered structure

Advantage of layered structure for file-system implementation is that it minimizes duplication of code. Multiple file systems can use I/O control and sometimes the basic file-system code. Each file system can then have its own logical file-system and file organization modules.

Limitations

The limitations of layering is that it can bring in more operating system overhead, which may lead to decreased performance.

Syllabus Topic : Implementing File System

5.2.2 Implementing File System

[Q] Explain the implementation of file system in detail.

Operating System (MtpSet)

5.2.2(A) File System Layout

Disk surface can be divided in many partitions. Sector 0 of the disk is called MBR (Master Boot Record) used to boot the system.

Each partition can hold the independent file system. Beginning and ending address of each partition is given in partition table which resides at the end of MBR.

When system is booted, it reads the active partition.

MBR- The MBR program reads in its first block, called the boot block, and executes it. Fig. 5.2.2 shows the possible file system layout.

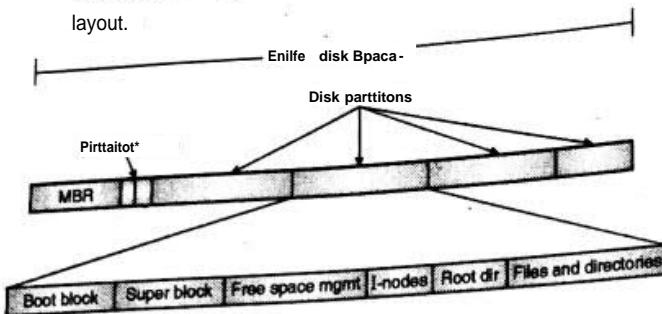


Fig. 5.2.2 : Possible file system layout

Bootblock

For homogeneity, every partition begins with a boot block, even if it does not contain a bootable operating system.

In future operating system might be loaded in that block.

Superblock

When system boots, all key parameters regarding file system is read in memory from this Superblock.

Usual information in the superblock includes a magic number to identify the file system type, the number of blocks in the file system, and other key administrative information.

Free space management

To keep track of free disk space, the system maintains a free-space list which records all free blocks.

I-node

The information regarding each file in file system is kept in data structure called I-node. For each file there is one i-node.

It is the top level file system

Files and directories

Files and directories in disk

5.2.2(B) Virtual File System

Write notes on virtual file system.

Q-

All modern operating systems must simultaneously support multiple file systems. The straightforward method of implementing multiple file systems is to use a virtual file system interface. This interface is the file system interface. This key layers- The first layer is the file system interface. This layer is implemented by descriptors. The second layer is called the virtual file system (VFS) layer having following two functions file system-generic operations from their

* Implementation by defining a clean VFS interface.

2. The VFS offers a means for uniquely representing a file throughout a network. The VFS is based on

which is a file-representation structure. The vnode denotes a numerical designator for a network-wide unique file. This network-wide uniqueness is needed for support of networked systems. The kernel keeps one vnode structure for each active node (file or directory).

Thus, the VFS differentiates local files from remote ones. The local files are again distinguished as per their file-system types. The activation of file-system-specific operations is carried out by VFS to handle local requests consistent with their file-system types and even calls the NFS protocol procedures for remote requests. File handles are built from the relevant vnodes and are passed as arguments to these procedures. The third layer of the architecture is the layer implementing the file system type or the remote-file-system protocol.

3. The object-oriented principles are used to design VFS. It has two modules: a set of definitions that states what file-system objects are permissible to be and software layer for these objects manipulation.

Following four major object types are defined by VFS:

i- An individual file is represented by i-node object

2. An open file is represented by file object

Flit	Attribute*
A	attributes
B	attributes
C	attributes
D	attributes

Fig. 5.2.3

- Fig. 5.2.3 shows simple directory containing fixed-size entries with the disk addresses and attributes in the directory entry.
- Some systems use entries are just name of the file and i-node number. All the entries in a file are stored in a single i-node.
- Fig. 5.2.4 shows this design.

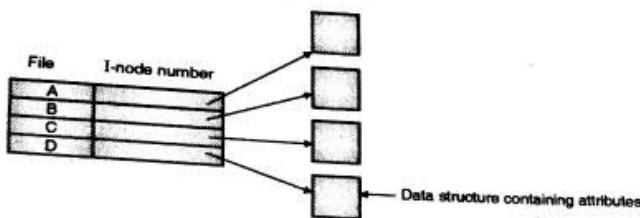


Fig. 5.2.4 : A directory in which each entry just refers to an i-node

- Above two approaches correspond to MS-DOS which have short file names of 1-8 characters and **optional extension of 1-3 characters**. In UNIX Version 7¹ file names were 1-14 characters, including any extensions.
- All modern operating systems support longer and variable-length file names. Following are the two approaches to implement it.
- The simplest solution is to keep file length of **maximum 255 characters**. The above two designs then can be used in which maximum 255 characters is **reserved for file name**. Since it is not common to have all files such longer names, directory space would be wasted.
- In this approach, fixed size of all directory entries is **not considered**. Each directory entry contains a fixed section, which begins with length of the entry, and then followed by data with a fixed format which are generally file attributes. It is shown in Fig. 5.2.5.

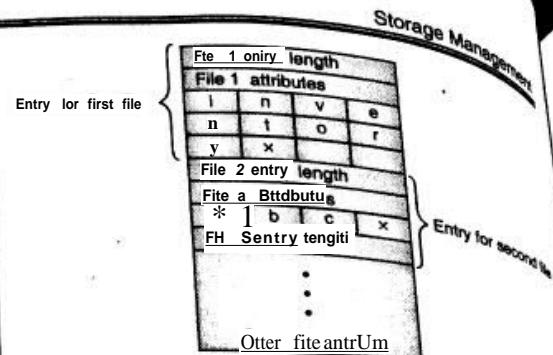


Fig. 5.2.5: In-line handling of long file names

- In this example two file inventory and abc is shown. Termination of each file is shown in Fig. 5.2.5 by cross symbol.
- A file entry in a variable-sized free space is initiated into the directory. It may happen that, next file to be accommodated in this free space may not be accommodated in this free space.
- Solution on this issue is compaction. Another problem is that a single file spreads in several pages, page fault reading of a file.

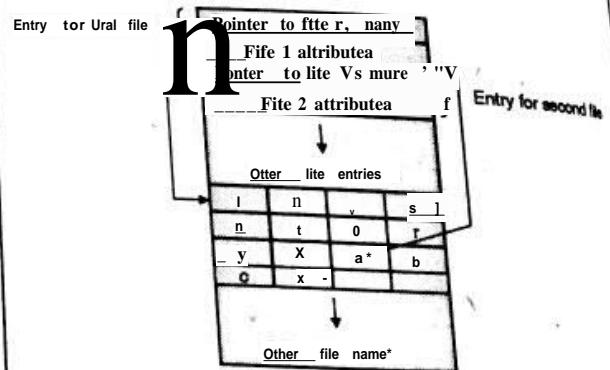


Fig. 5.2.6 : Long name handling using heap

- The second approach to deal with variable-length names is to allow directory entries of fixed length and maintain the file names jointly in a heap at the end of the directory as shown in Fig. 5.2.6.
- The drawback discussed above is defeated here. When an entry is removed, the next file entered will fit in freed space. The management of heap is burden here and page faults cannot be avoided as in previous method.
- All the designs discussed, suppose searching directory sequentially from beginning to end.

tiding up the file name. Search is used if directories are found in the linear list memory. jamming the linear list memory. more time is used to create a new file, searching for the same name. A file, that is released. of the directory file is released. Following many options can be used.

- o The file entry can be marked as unused.
- o Assign a special name to fit, blank name.
- o A used -unused bit in each entry can be maintained.
- o **Attach it to a Parent directory**

Adoption is to copy the last entry into the released location and to redice the length of the directory. A linked list can also be used to decrease the time required to delete a file.

Solution to this problem is use of hash table in each directory as hash search is efficient than sequential search.

5.2.3(B) Hash Table

In this method, a linear list is used to store directory entries. In addition to this, a hash data structure is used. The hash table uses a value computed from the filename and returns a pointer to the file name in the linear list. As a result, it can significantly decrease the directory search time. Insertion and deletion are also quite simple. It is required to make some provision for collisions-situations if two file names hash to the same location.

The main problems with a hash table are its usually fixed size. The hash function is dependent on size of hash table. Consider a linear-probing hash table which holds 64 entries. Using hash function, file names are converted into integers from 0 to 63. For creating a file, we must increase the size of directory hash table for example to 128 entries. Consequently, we need a new hash function which will convert file names to the integers ranging from 0 to 127. In this case, it is necessary to reorganize the present directory to reflect their new hash-function values.

Alternate solution is to use a chained table. In this table, a linked list is maintained instead of an individual value.

be handled by adding a new file, searching for the same name. A file, that is released.

Although the file is released, the link is still present in the directory. This is known as incomplete directory.

5.2.4 Allocation Methods

Q. What are the allocation methods? (Dec. 14, May 16)

MU - Dec. 2014, May 2016, 10 Marks

Disks supports for direct access, which we can have flexibility in the organization of files. There should be space to fit the files. The space should always allow for the access of the files. So that quick methods are used majorly in different operating systems.

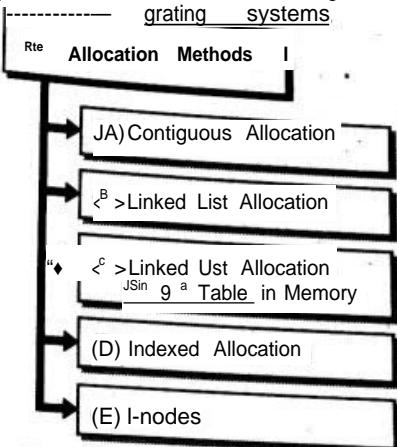


Fig. C5.9 : Allocation methods

5.2.4(A) Contiguous Allocation

Q. Explain contiguous file allocation method with its advantages and disadvantages.

- In contiguous allocation each file takes up a set of contiguous blocks on the disk. Disk addresses define a linear ordering on the disk. If each block size on disk is 2 KB, then a 30 KB file would be allocated 15 consecutive blocks.
- Contiguous allocation of a file is defined by the disk address of the first block and length (total blocks occupied). If the file is m blocks long and starts at location i , then it takes up blocks $i, i + 1, i + 2, \dots, i + m - 1$.

The directory entry for each file specifies the address of the starting block and the total number of blocks allocated for this file.



- Directory entry in Fig. 5.2.7 shows that file A starts at block 0 and it is 3 block long occupying block 0, block J and block 2. Similarly, file B starts at block 6 and it is 5 block long and so on.

Advantages

- It is simple to implement because, only information needed to keep track on files block is starting block and length (total number of blocks occupied).
- For sequential access, the file system memorizes the disk address of the last block referenced and, when required, reads the next block.
- To access block k of a file directly which starts at block i , we can immediately access block $i + k$. Thus, both sequential and direct access can be supported by contiguous allocation.
- Performance is high as whole file is read from disk in single disk operation. To reach to the first block only one seek is needed.

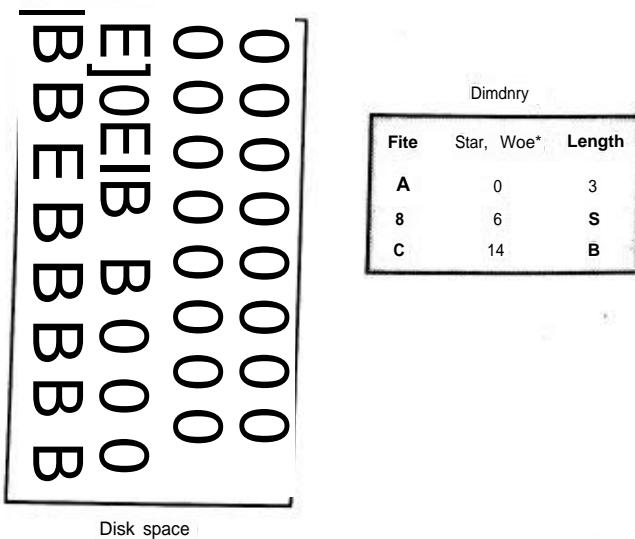


Fig. 5.2.7 : Contiguous Allocation of disk space

Disadvantages

- When allocated file is deleted, continuously allocated blocks to the file become free. For new file to allocate, these blocks might not be sufficient. If new file size is small than previously allocated size, some holes remains which are not enough to allocate any new file. As a result, the disk ultimately consists of files and holes causing external fragmentation.
- In some cases, it is simple to find out space is needed for a file to be allocated. For output file this size estimation becomes difficult. For CD-ROM, contiguous allocation is reasonable as all file sizes are known in advance.

5.2.4(B) Linked List Allocation

- Q. Explain linked list file allocation method with its advantages and disadvantages.

- **Linked list allocation overcomes the limitation of contiguous allocation.** In linked allocation, each block of a linked list of disk blocks. The scattered disk blocks can be allocated to the file. The directory entry holds a pointer to the first and last blocks of the file. Creation of a new file is simple. For this, it is required to create a new entry in the directory. The directory entry holds a pointer to the first disk block of the file. The pointer value is nil for empty file and file size for this file is zero.
- The pointers need to be followed from block to block in order to read a file.

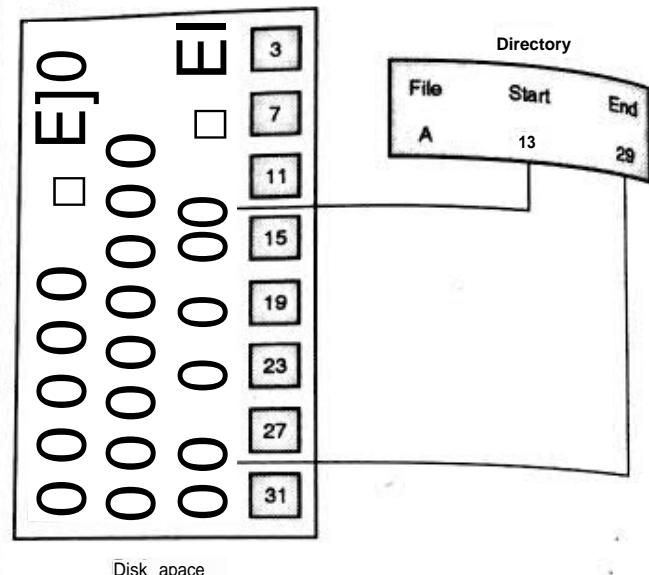


Fig. 5.2.8: Linked List Allocation of disk space

- Fig. 5.2.8 shows linked list allocation for file A. The file A of five blocks start at block 13 and continue at block 16, then block 21, then block 24, and finally block 29. Each block holds a pointer to the next block. These pointers are unavailable to the user.
- User can see the block size excluding the size required to store the address. Accordingly, if each block size is 512 bytes, 4 bytes are required to store address, then the user sees blocks of 508 bytes.

Advantages

- In linked list allocation, every disk block can be used which is not possible in contiguous allocation. Hence there is no external fragmentation except for internal fragmentation in the last block.
- Reading a file sequentially is simple.
- It is not necessary to compact disk space.

Z J, dva^{nta} 9^{es} 5-2

Z o Z ineff^{cient} to su Pport a ~~part~~ m acce ss c*
retrieve to block i_t the onemf^g at the beginning and read the i - 1 blocks previous to it,
To one at a time.

In each block pointer takes some space, so each file
require slightly more disk space rather than it[®] actual
j j e.

if pointer is lost or damaged f_tJe
essibJe
WouId become

jj ,

A tfromg pointer can be chosen due
operating ^{evstem} ^{Software} system breaka^o bUg « the
hardware. As a result of ww ch. ^{br}wr, of disc
result in unkng into the ftee- space list
file, could tun, into another

5.2.4(C) Linked List Allocation using a Table in Memory

In linked list allocation, each block ~~contains~~ information, therefore entire bkt used to store file content. This linX.J! ~~is~~ ^{now} fully eliminated by keeping pointer inform ~~ation~~ ^{ation} in bch always remains in memory, ⁱⁿ Ube

If file A occupies block 5, 9, 1, 12 and 7 in the same order and file B occupies block 3, 8, 4, 14 and 6 in the same order. The File Allocation Table (FAT) is shown in Fig-5.2.9.

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

Fig. 5-2.9: Linked List Aitoa**¹⁰⁰

Advantages

Storage Management

Disadvantages

work. Ubie TM USt ** *n “*TM»y all the time to make it .

5.2.4(D)

Explain indexed file allocation and disadvantages. With advantages and disadvantages.

list allocf AIIOCatiCITI Table AT) in memory, linked table m o!l SUPPORT raid0Tn access , but the ent ire allocation. all the pointers all the time. In indexed called as index block.

The : a **jiuca djock** assigned to each file and this mdex block holds the disk block addresses of that particular file. There is a pointer from i^{th} entry in the index block to the 1^{th} block of the file, it means n^{th} entry in index block holds the address of the 0^{th} disk **block**. **The address of the index block of the file is** maintained in directory.

In order to locate and read the i^{th} block, pointer in the i index-block entry is used.

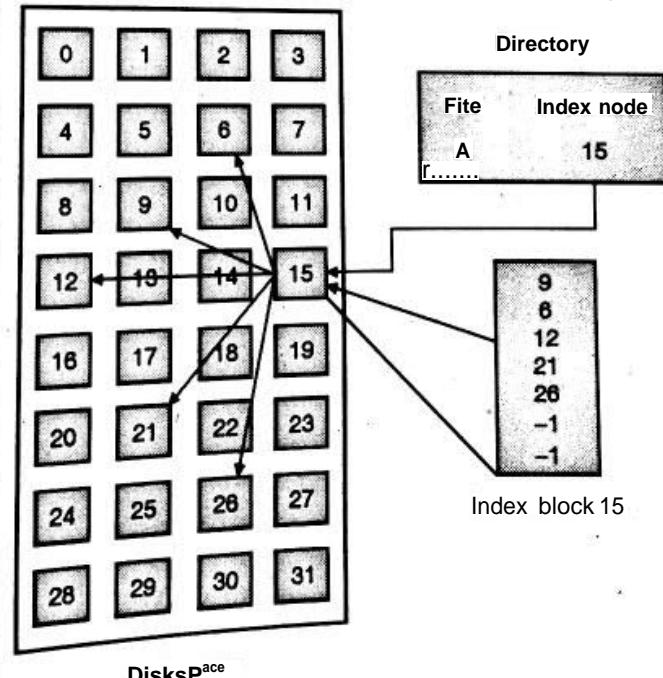


Fig 5.1W: Indexed Allocation of disk space

Initially when the file is created, all pointers in the index block are initialized to nil. When the ith block is first written, a block is obtained from the free-space manager. and its address is put in the ith index-block entry.

Advantages

- Indexed allocation supports random access.
- There is no external fragmentation, because any free block on the disk can be allocated as per request for more space.

Disadvantages

- The pointer overhead of index block is more with compare to the pointer overhead of linked allocation.
- Indexed-allocation schemes suffer from some of the same performance problems as does linked allocation.

UNIX Example

The structure of i-node

Q- Explain the structure of i-node in detail.

The information regarding each file in file system is kept in data structure called i-node. For each file there is one **i-node**. An active i-node is associated with exactly one file, although many file names may be associated with a single i-node. Every file is prohibited by precisely one i-node.

Operating system manages all types of file by means of i-nodes. The exact i-node structure differs from one implementation of UNIX to another. The file attributes as well as its access rights (permissions) and other control information are stored in the i-node. The FreeBSD i-node

Storage Management
structure is shown in following Fig. 5.2.11. It contains the following data elements:

- The mode of access and type of the file.
- Identifiers needed for group access and h₀ i₁ file.
- The time at which the i-node is just now system and the creation time of the file.
- The time at which file is just now read and written.
- The file size which is in bytes.
- The series of block pointers.
- Total number of actual physical blocks Octaiae file. It also comprises the blocks holding pointers and attributes.
- Total number of directory entries which refers file.
- The flags which user or kernel can demonstrate the uniqueness of the file.
- An arbitrarily chosen number allocated to every time that the latter is assigned to a new fit as generation number of the file. The same number used to identify references to deleted files.
- The data blocks size pointed by the i-node. Usually it is the same as, but sometimes larger than, the file system block size. FreeBSD has a minimum block of 4,0% bytes (4 Kbytes).
- Extended attribute information size.
- Zero or more extended attribute entries.

I-node table, or i-node list present on disk, contains i-nodes of all the files in the file system. When a file is opened, its i-node is carried into main memory and stored in a memory-resident i-node table.

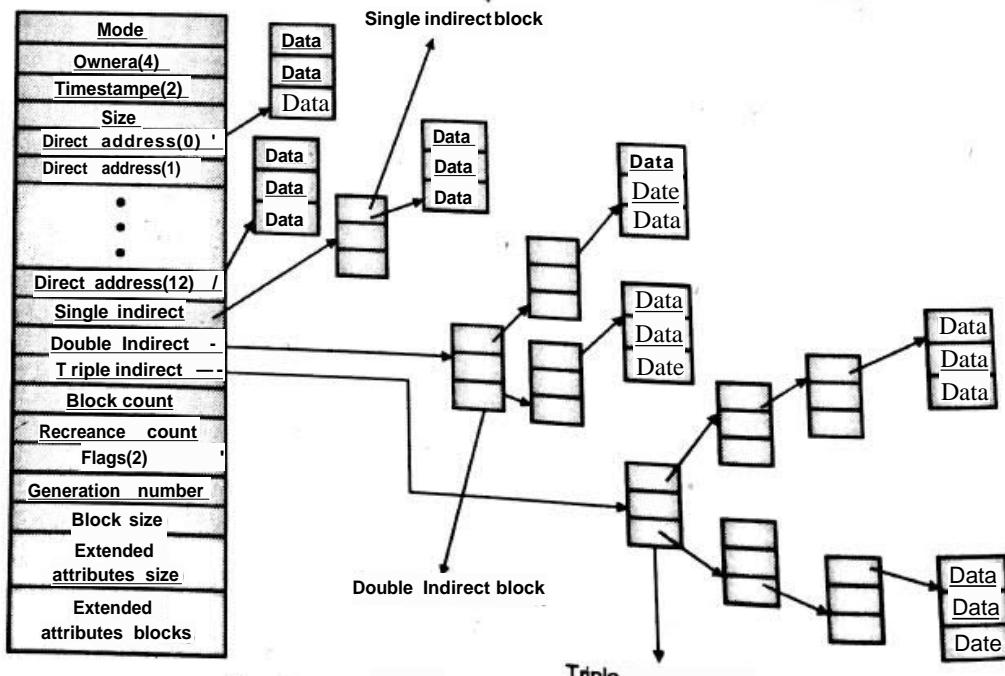


Fig. 5.2.11 FreeBSD i-node and file struct.

5.2.4(E) /nodes

Q. How i-node is used to allocate the file?

I-node (inode-node) is used to store the file address, file name, file type, file size, and file modification time. It also contains pointers to the disk blocks that store the file data. The i-node is located in the file allocation table (FAT) in memory. When a file is opened, its i-node is loaded into memory, and the file's data is read from the disk blocks pointed to by the i-node.

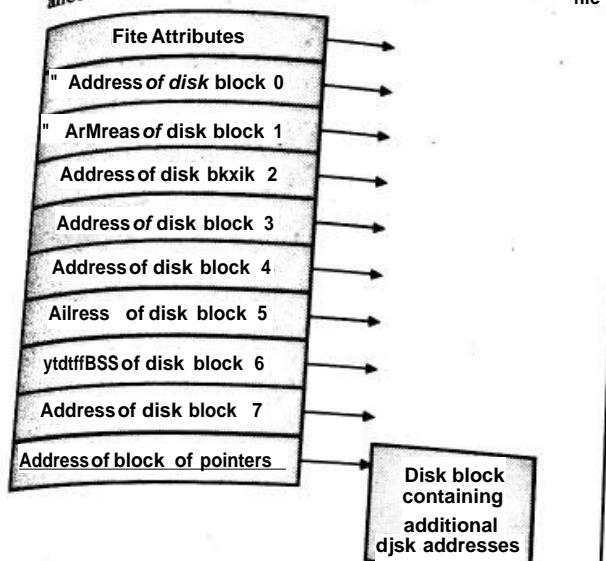


Fig. 5.2.12 : Example of i-node

- the size of the file is proportional to the number of blocks that disk contains. If disk has k blocks, file allocation table will have k entries,
- in contrast, the i-node scheme requires an array in memory whose size is proportional to the maximum number of files that may be open at once.

Fig. 5.2.12 shows example of i-node.

- the last disk address is reserved for the address of a block containing more disk block addresses.

The reason is, each i-node has space for fixed number

of disk addresses. If file grows beyond this limit, this list of disk addresses will point to disk blocks containing additional disk addresses.

Topic : Free Space Management

Free Space Management

This

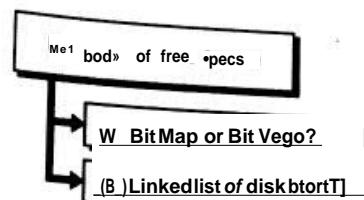
When a file is deleted, the disk space becomes free. This space should be reused to allocate to new files.

keeps track of free space. Store Mana merit

blocks *hi ch '77 usm @*** list, The disk free. File sp elis t records 2 To .***-e a neJ r , UCh bl ks availability o, foM sPacc lira 1. searched for gets located to Z TM Spce If fooad. free space file sp ** Im F** W File " d " <>* from the again gets, cl ud ?r SPS " due 10 " deletion of file is ** W!. Q th -- space management

Q. List methods of free space management.

To implement free space management, following two methods are used:



Fc- C5.10 • Methods of free space

** 5.2.5(A) Bit Map or Bit Vector

Q. Explain bit map method in detail.

- In this method, every block on disk is signified by a 1 bit. A block is represented by bit 1 and an allocated block is represented by bit 0.
 - Suppose block numbers 3, 5, 7, 9, 10, 15, 18, 20, 23, 25, and 28 are allocated and rest of the blocks are free. The free space map would be as follows. Allocated blocks are shown in bold space.
- 11101010100111101101011010110.....
- With a bitmap, it is also possible to keep just one block in memory, going to disk for another one only when it becomes full or empty. Since the bitmap is a fixed-size data structure, if the kernel is (partially) paged, the bitmap can be put in virtual memory and have pages of it paged in as needed.

Advantages

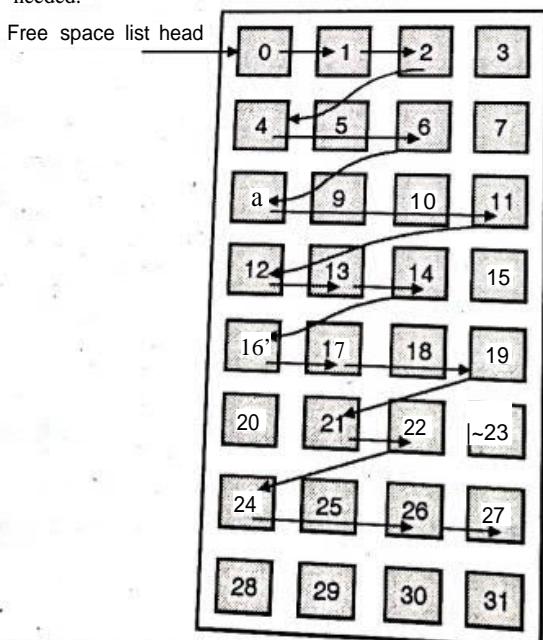
- It is relatively simple and requires less space, since it uses 1 bit per block.
- It is efficient in searching the first free block or k successive free blocks on the disk.

Disadvantages

- The whole vector is needed to keep written to disk. Otherwise it would be inefficient occasionally for recovery needs.
- Keeping bit map in main memory for larger smaller disks. It is not necessarily proportional to disk size.

5.2.5(B) Linked List of Disk**q. Describe linked list of disk block method**

- In this method of free space management disk blocks are linked together.
- The pointer to first free block is kept in part location on disk and cached it in main memory first block contains a pointer to the next free block, and so on.
- In our bit map example, we would keep a pointer to first free block which is block 0. Block 0 would contain the pointer to block 1, which would point to next free block, which is block 2 and so on. In Fig- 5.2.1 shown, block 0, 1, 2, 4, 6, 8, 11, 12, 13, M, 1b. 17*1 > 21, 22, 24, 26, and 27 are free blocks.
- This method is not efficient because, we must read each block in order to traverse the list, which requires extensive I/O time. But advantage is that, traversing the free list is not regular operation.
- Usually, the operating system simply needs a free block so that it can allocate that block to a file, so the first block in the free list is used.
- The FAT method includes free-block accounting into the allocation data structure. No separate method is needed.

**Fig. 5.2.13 : Linked free-space list on disk****5.2.5(C) Grouping**

This method is a modification of the free-list method. In this method, the first $n-1$ of these blocks are free. The first block stores the addresses of the free blocks. The addresses of a large number of free blocks can be stored in a linked list.

5.2.5(D) Counting

Allocation algorithm is used for space allocation or freed simultaneously. Hence, instead of maintaining a list of n free disk addresses, it is advantageous to keep the address of the first free block and the number (n) of free contiguous blocks. The entry in free-space list denotes a disk address and a count. In this approach, more space is needed for each entry with compare to a simple disk address. But advantage is that, the overall list is shorter, as long as the count is generally greater than 1.

Instead of linked list, a B-tree can be used to store these entries so that, efficient lookup, insertion, and deletion can be carried out.

Syllabus Topic : Efficiency and Performance**5.2.6 Efficiency and Performance**

Q. Explain various techniques to improve efficiency and performance of secondary storage.

As disk is slowest memory component, it may become bottleneck in system performance. There are various techniques to improve efficiency and performance of secondary storage.

5.2.6(A) Efficiency

The disk allocation and directory algorithms mainly decides the efficient use of disk space. The preallocation of UNIX inodes is done on a volume. This results some space occupied by inodes although disk is empty. Preallocating the inodes and scattering them across the volume improves file system performance. The UNIX allocation and deallocation algorithms maintains a file's data blocks near that inode block to lessen the seek time.

It is necessary to consider the types of data kept in a file's directory (or inode) entry. Normally, the records "write date" gives the information that is user to decide whether back-

done or not. Some system "P of file" <ot date." Which is useful to use the last access as last read.

As this information is maintain

writing the file in the directory > Astern. j(k- file * Block " essential to *henev* then if section Changes, and if the file be back out to disk, as place only in block chunks. CX d-n S On disks ta situation, if any time a file is o " practical its directory entry must be read and wri £ " "8 <n

s necessity cannot be efficient for as often accessed frequently. Therefore, we must benefit against its performance cost when designing a file system. In general, all file related data items needs to be considered for its effect on efficiency and

The size of pointers used by systems is different. If system uses 16 bit then length of file restricts to 64 KB. For 32 bit pointers, file length supported restricted to 4GB. Some systems use large size pointers, which occupies more disk space.

5.2.6(B) Performance

Accessing the memory is faster with compare to accessing the disk. For accessing single word, the memory access is on the order of a million times as fast as disk access. Due to this differentiation in access time, many file systems have been designed with various following optimizations to improve performance.

Various optimizations to Improve performance

- 1. Caching
- 2. Block Read Ahead
- 3 Reducing Disk Arm Motion

Fig. C5.11 : Various optimizations to improve performance

1. QW " 8
Block cache or buffer cache is the most 8
technique used to minimize disk accesses. In

Storage Management

101ht d&k w,* 8rou Porbi, cles that logically belong to Perth " e *,>< Cpt in mory for the and utilized to handle the c *Clic, *8orithms exist The most 8w ,41 algorithm is to make certain about all c*he?* SIs to notice if the needed block is in the fulfill. If it is in cache, the read request If lbc b| * k, Ut a disk access.

cache, "d then coXd t U " fTM >< the >, sqoes sive r Xs f T" " " needed Aft " h,lr,U 'd f ro >n the X'e In cache there

,KSKS ar y 10 h a v? n be lar8e " umber o f bl <ks. It is 1 bl <k. Hash , hl searching technique to locate blocks >1 cache , orga in @ is used to organize value are linked together

Elision n chain form g a linked list so the can be followed,

To bring a new block

needed 1010 1 fUU Cache , some block " been m r rem o Vcd iU,d transfer to to disk , If U has been SlKC being brou 8 ht in TMs circumst . is very much Hke P3 @ 118 * 411(1 all the normal page replacement algorithms such as FIFO.

LRU Wopnate to replace block from cache.

Since cache references are relatively occasional, it is practicable to keep all the blocks in exact LRU order with linked lists. Fig. 5.2.14 shows buffer cache data structures.

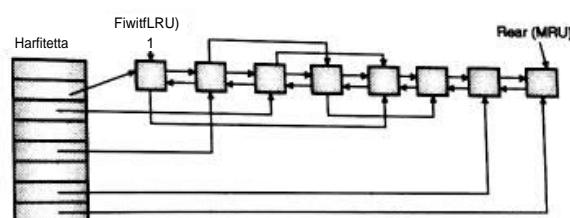


Fig. 5.2.14 : The buffer cache data structures

The Fig. 5.2.14 shows a bidirectional list running through all the blocks in the order of usage, with the least recently used block on the front of this list and the most recently used block at the end of this list. When a block is referenced, it can be removed from its position

on the bidirectional list and put at the end. In this way, LRU order can be maintained. It seems that LRU is Arable, treasons are;

i-node like critical block is read into the cache. If it is modified but not rewritten to the disk, a crash will leave the nlesystem in an inconsistent state. H the i-



- node block put at the end of the LRU chain, it may be quite a while before it reaches the front and is rewritten to the disk.
- Some blocks, such as i-node blocks, are not often referenced two times within a small interval. In order to consider the above two factors regarding crashes and file system consistency, modified LRU scheme is possible. The blocks can be categorized as i-node blocks, indirect blocks, directory blocks, full data blocks, and partially full data blocks. Blocks that will most likely not be required again quickly go on the front, instead of the rear of the LRU list, so their **buffers will be reused quickly**.
- Blocks that might be required again soon, such as a partly full block that is being written, go on the end of the list, so they will stay around for a long time. If the block is necessary to the file system consistency and it **has been updated, it should be written to disk straight away, in spite of which end of the LRU list it is put on.**

→ 2. Block Read Ahead

- The performance of file system can be improved by **bringing block in cache before it is required**. Whenever file system creates n^* block, it checks if $n+1^*$ is **available in cache or not**. If not available, it brings it in cache taking into consideration that, it would be required in future.
- This read ahead policy does not work for random access and only works for files that are being read sequentially. The read block in cache would force to replace the useful block and it might not be required. The file system can keep track of the access patterns to each open file. If next time file read is suppose sequential and not random, the read ahead strategy is useful otherwise not.

→ 3. Reducing Disk Arm Motion

- Apart from the caching and block read ahead technique to improve performance, other technique is to minimize the total disk arm motion by keeping blocks that are likely to be accessed in succession close to each other, if possible in the same cylinder. In the following ways performance can be improved
- If bitmap technique is used to record free blocks, and the entire bitmap is in main memory, then it is simple to select a free block as close as possible to the previous block
- Grouping of the blocks can be done with free list. Now, system will keep track of disk storage in groups

of consecutive blocks instead of only block. sedor holds 512 bytes in the sectors. coul? / ot, blocks occupying 1 sectors but assign disk storage units of 2 blocks occupying 4 sectors. In this case, cache would still use 1-KB blocks, disk transfers would still be 1 KB but reading sequentially would minimize the number of seek. The factor of two, considerably improving performance by variation on the same theme is to take Kc J. A rotational positioning. When allocating blocks, system tries to place consecutive blocks in a file in the same cylinder.

Those systems which use i-node like data structures, that two disk accesses are required to access a short file: one for i-node and other for data. Generally i-nodes reside at the start of disk having average distance from its blocks about half the number of cylinders.

Because of this reason long seek is required. Performance can be improved by placing i-nodes in the middle of disk which would minimize the average seek between the i-node and the first block by a factor of two.

Other approach to improve the performance is to divide the disk into cylinder groups, each with its own nodes, blocks, and free list. When creating a new file any i-node can be selected, but a try is made to search a block in the same cylinder group as the i-node. If nothing is available, then a block in a nearby cylinder group is used.

Syllabus Topic : Recovery

5.2.7 Recovery

- Q. Explain file system recovery in detail.

Directories and Files reside on disk or remains in main memory. System crash may lead in data loss or in data inconsistency. Following are the ways to handle these issues.

Ways to handle Recovery Issues

- (A) Consistency Checking
- (B) Backup and Restore
- (C) log Structured File Systems

* Ifr C5.12 : Ways to handle recovery issue*

In order to get the per 06 ""9 information, it is cache in ~~acc~~ ⁰⁶ updated copy of directory information ~~in~~ ^{*-Pt in T} ^{in k di}, ~~ect~~ ^{memory. Whereas correspon} ^{memory. g.} ^{resides in main} ^{information} ^{o, the} ~~ssen~~ ^{<**} is old copy as cached ~~disk~~ ^{dis} ^{information} ⁿ ^{o, the} ^{tially written to disk on} ^{-y info} ^{o, the} there can be loss of co ~~nta~~ ^{nta} ^w ^{Update} ^{is not} also of TO operations ^{curn} ^{J.} ^{Cacl} ^{*} ^{and h} ^{*} ^{anq} ^{ncbinc} ^{crash.} As a res ^{ult} ^{n, y} ^{Wor} ^{er. anq} ^{rones} ^{of opened f} ^{iles} ^y ^{Uion, f} ^{af} ["] ^a ^{adi} ^{the file system i} ⁿ ^M ⁻ ^{** lost Th} ^{in the} ^{state the actual state of son*} ^{**} ⁱⁿ ^{nsistent} ^s [,] ^{Won} ^{compare to present i} ⁿ ^{t.} ^{***} ^{'s differ} ⁱⁿ ["] ^{His} ^{regularly, during reboot specie} ^{rec} ^{try} ^{,"} ^{*hh} ^e ^{and correct disk 'neons, f"} ^{gram verify} ^{tfe} ^{consistency checker p r o g, f*} ^{The} ^{kdsk in MS-DOS.} ^{8 dms} ^{ft/sck} ⁱⁿ ^{»f}

The directory structure data and
are compared b * the > prtg rams , ? 10 Cks Of > disk
KCU rs then tries to correct any incor any on any match
The disk space allocation and free onsj WK stencie if exist
algorithms direct the checker proe/"/"! 7 ! element
type of problems and how - to dis CO' er the
correcting them. SUCC - M it will be in

in case of linked allocation, a link is block to its next block. ThereforT ^{6_*"} My reconstruct the file from the data blwfc ^{is easy} to directory structure is created again After this ,

In case of an indexed allocation, loss of a block is difficult to handle. This is because the blocks are unknown about the other data blocks. Therefore, for read accesses, UNIX caches directory entries. On the contrary, for write leading to space allocation, or other metadata updating, is carried out synchronously, prior to the corresponding data blocks are written.

*5.27(8) Backup and Restore

- If failure of magnetic disks occurs then loss of data should not be forever. Back up of data from disk to storage devices can be carried out by system programs. Recovery of any data is now just the restoring the data from backup.

The amount of copying can be minimized by making use of information from each file's directory entry. I k** up program is aware of the time of the last bac. P of a file, then file's last write date in the ~~current~~ date, then denotes that the file is unchanged since that there " no need to copy the file again.

le Systems

Applied effective recovery algorithms can be 'becking. The im_D[eni, , , problem of consistency 2TM «UTMor] The disk file (nLd <TMT' resulted in 1₀ KbaMd syste jounalling) fil e systems. dit «'y structures fr e,7? SmKturus for «>pl K **Pointers** can leads' £ T* **Pointers, and** f₀ FCB failure of system. "consistent state due to **Changes**

place whenTise" 2?* 10 ab o Ve 5tn,ctures ■»
incorporated in n~. « 8 „based techniques was not
results in 0PeCab " S systems ■ The creation of file

X on^{ra} Z S X Ur te " Wi,hin,heBle
modificaon of dhereotry struck, XoXrTs

• BLOCKS and therefore free counts for all these blocks are decreased. If crash occurs then the changes can be interrupted and structures can lead to inconsistent state. These structures can be fixed after recovery. But in such scheme following difficulties may arise .

The inconsistency may be beyond the repair. The consistency check made may be incapable to recover the structures, which can lead in loss of files and yet whole directories. In order to get rid of the conflicts, a human involvement is required. Availability of system will not remain to the user until the human directs it about how to proceed. For the consistency checking of huge size data can require hours of clock time to check.

If log-based recovery techniques are applied to the system, metadata updates then the above problems can be solved. NTFS and the Veritas file system use

- the log-based recovery techniques All changes done in **metadata** are written sequentially to a log.
 - These written changes to the log considered to be committed. After this the control of execution of system call can return to the user process which then carries on execution. In the meantime, replaying of log entries are carried out across the actual file system structures. Because the changes are done, a pointer is updated to show which actions have finished and which are still unfinished.
 - Once transaction finishes the execution and committed **entirely**, it is detached from the log file. This log file is nothing but a circular buffer. In circular buffer writing takes place from start to end and when reaches at other end it then continues at the beginning, overwriting older values as it goes.
 - After the crash of the system, there can be no or any number of transactions in log file. If the transactions present currently in log file were not completed to the file system, although operating system has committed it, they must now be completed. The transactions can begin execution from the pointer until the work is finished so that the file-system structures stays in consistent state.
 - If transaction does not committed prior to system crash and aborted then Any changes that were performed to the file system must be undone. So the file system will remain in original state and consistency will be maintained. The logging on disk metadata updates are **performed faster than with compare to applied straight** to the on-disk data structures.

Syllabus Topic : NFS

5.2.8 NFS

NFS is developed by Sun Microsystems and it is used on Linux to join file systems of different computers into one logical whole. Version 3 of NFS was introduced in 1994 NSFV4 war introduced in 2000 and offers a number of improvements over the previous NFS architecture.

5.2.8(A) NFS Architecture

- Q. Explain operation of netwkWesystem in detail.

- NFS allows a number of clients and servers to share a common file system. These clients and servers can be located geographically at long distance. Every machine to be both a client and a server simultaneously.

Remote clients access the NFS network. Every NFS server exports directories. Every NFS server exports more than one directory.

Stora ge Management

its **directories**. These **directories** then **acCeiv**
remow clients. **directories** along **sed by**
subdirectories, so in fact entire **directory** **U**
usually exported as a unit. **J** **—** *** —** **te**

The directories list that server exports is kept in a file, often */etc/exports*. As a result of this, these directories are exported automatically whenever the server starts.

ed. Clients can access exported directory by them. When a client mounts a remote directory, its directory mounted by Xw NFS **. It becomes part of the hierarchy and client can access these directory.

5.2.8(B) NFS Protocols

- As in the network, client and server may have different operating systems on different hardware, environment is heterogeneous. It becomes necessary that interface between the clients and servers be defined.
 - Then and then it is possible for everyone to write a new client implementation and guarantee it to work correctly with existing servers, and vice versa. Up to achieve this objective by defining two client protocols. The mounting is handled by first NFS protocol.
 - A client can send a path name to a server and request permission to mount that directory somewhere in the directory hierarchy. The location at which it is to be mounted is not specified in message from client to server.
 - The reason is, servers do not care for this mounting. If the path name is legal and the directory specified has been exported, the server returns a file handle to the client. The file handle includes fields exclusively identifying the file-system type, the disk the i-node number of the directory, and security information.
 - Succeeding calls to read and write operations in the mounted directory or any of its subdirectories use the file handle. i.e. directory replication is not supported by NFS, so user may arrange for all the file systems to be the same.
 - As a result, automatic mounting is preferred for read-only file systems with system binaries and other files that hardly ever change.
 - Example:
of nfc, an access protocol is used. Clients can send purpose messages to servers to manipulate directories and read and write files. The also access Z time, file attributes, such as file mode, size, UNIX I as max file size protection mechanism is also used in NFS.

5.2.8(C) NFS Implementation

5-31

a- Explain the layered structure of NFS.

Fig. 5.2.15 shows NFS layered structure. The system-call layer is the system-call layer in which file system (VFS) layer is the system-call layer. The VFS layer maintains a mapping between file names and file handles. The NFS client calls the NFS server to read or write data. The NFS server returns the data to the client.

The VFS layer maintains a mapping between file names and file handles. The NFS client calls the NFS server to read or write data. The NFS server returns the data to the client. The NFS client maintains a buffer cache to store data. The NFS server also maintains a buffer cache. The NFS client sends messages to the NFS server and receives messages from the NFS server. The NFS server sends messages to the NFS client and receives messages from the NFS client. The NFS client and NFS server both have local file systems.

Consider system calls like `mount` as an example. The `mount` program is called by the administrator for a directory, the local directory on which the remote directory is mounted, and other information.

The `mount` program parses the directory to be mounted and finds out the NFS server on which the remote directory is located. It then contacts that machine, asking for a file handle for the remote directory. If the directory exists and is available for mounting, the server returns a file handle for the directory. Finally, it makes a `mount` system call, passing the handle to the kernel.

The kernel then constructs a v-node for the remote directory and asks the NFS client code to create an r-node (remote i-node) in its internal tables to hold the file handle. The v-node points to the mode.

Each v-node in the VFS layer will ultimately contain either a pointer to an r-node in the NFS client code, or a pointer to an i-node in one of the local file systems. Thus, from the v-node it is possible to see if a file or directory is local or remote.

If it is local, the correct file system and i-node can be located. If it is remote, the remote host and file handle can be located.

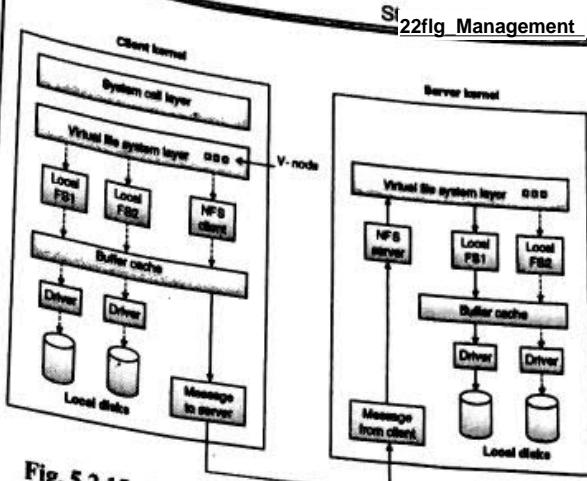


Fig. 5.2.15 - Layered structure of NFS

Syllabus Topic : Secondary Storage Structure

5.3 Secondary Storage Structure

The secondary

West level of the system. structures are the

Topic : Overview of Mass-Storage Structure

5.3.1 Overview of Mass-Storage Structure

This section

of tertiary storage devices

structure of secondary

5.3.1(A) Magnetic Disks

L-51 Ex plain physkai uare of magnetic disk

The magnetic disk is used as main storage device in your computer system. It is a magnetic type of storage device. Within one magnetic disk unit, many physical disks are present. Each disk is known as a platter. Several metal platters are present within the magnetic disk and these are coated with a special magnetic material. The platters rotate many thousands of times within a second for accessing the data. The diameters of platters are usually between 1.8 to 5.25 inches. Magnetic read and write head are present which floats just above the surface of the platter.

Material used for making the platters is Aluminum, glass, or ceramic and two read/write heads are present, one for upper and lower surface. Platters are arranged in stack because of which the position of the read/write heads often is referred to by its cylinder. Cylinder is

Operating System (MU - Sem 4 - IT)

the location of a single magnetic disk can be read from and write to disk cache mechanism. Now data is the part of magnetic disk. Magnetic disk can be used as separate disk for each partition. The spinning speed of most of the disk is not use. The rate at which data flow between disk and the computer is called as The transfer rate. The random access time, cylinder the time to move the disk arm to the required cylinder (seek time), and the time for the required sector to rotate to the disk head (rotational latency).

disk can have seek times and rotational latencies of several milliseconds. It is possible that the head may make contact with the disk surface. Sometimes it is also possible that head may damage the magnetic surface called head crash. After head crash, the entire disk must be replaced because it cannot be repaired.

A removable disk permits different magnetic heads to be mounted if needed. Removable disks contain one platter, covered in a plastic case to avoid damage while not in the disk drive. Floppy disks are inexpensive removable magnetic disks that are a soft

plastic case containing a flexible platter. The head of a floppy-disk drive generally sits directly on the disk surface, so the drive is designed to rotate more slowly than a hard-disk drive

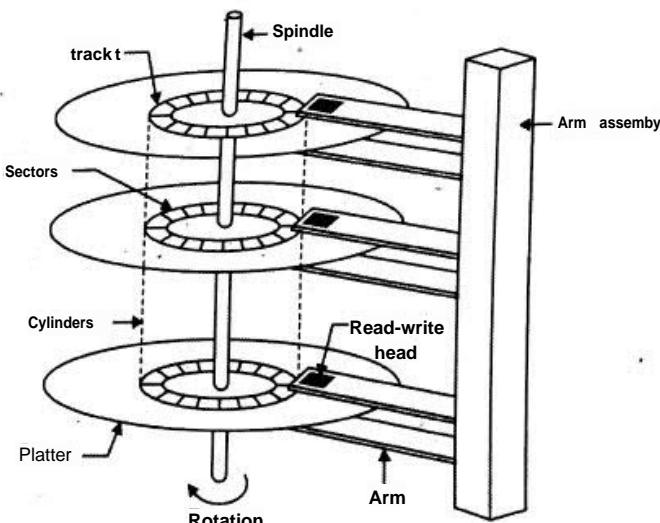


Fig.53.1: Moving-head disk mechanism

The storage capacity of floppy disk is only 1.44MB. It is a small plastic plate coated with magnetic material. Data is recorded in floppy disk in terms of magnetic spots. We can transfer the data from one computer to

other using parallel on floppy disk. Its size is 3 1/2 inch which is floppy is 3W inch. Hard plastic cover is used to cover it. A metal shield is used to close the read/write window. The disk is in a TX drive. A notch is used to write protection. A floppy disk has a size of 5 1/4 inch disk. The opacity of 1.2 mb. It is a 5 1/4 inch plastic plate coated with magnetic material.

5.3.1(B) MagneticTapes

Q. Explain physical structure of magnetic tape.

Magnetic tape is the example of an early secondary medium. 71tb of tape is comparatively slow and have large capacity, its access time is slow than that of RAM and magnetic disk. Magnetic tape is thousand times slower for random access compared to random access to disk. So disk is more preferred secondary storage with compare to tape. Primarily, the main use of tape is for taking backup and for storing the infrequently used information. It is used to transfer information from one system to another. A tape is kept in a spool and is wound or rewound past a read-write head.

For moving to the desired position, tape may take time in minutes. But once located, the write operation is carried out with speed equal to speed of disk drive. Tape capacities are decided by particular type of tape drive. Usually, tape storage capacity ranges between 20GB to 200GB. The storage capacity is double for some tapes which are having built-in compression.

Syllabus Topic : Disk Structure

5.3.2 Disk Structure

- Modern disk drives are viewed as large one-dimensional arrays of logical blocks. This logical block is the smallest unit of transfer. Logical block size is usually 512 bytes. The low-level formatted disks can have a different logical block size, for example 1,024 bytes.
- The mapping of one-dimensional array of logical blocks is carried out onto the sectors of the disk in sequential manner. On the outermost cylinder, in first

track, the very first sector is sector 0. The mapping is carried out in order through this track, then through the remaining tracks in that cylinder, and then through the remaining the cylinders from 0 to 15.

fW conversion of logical **address** into **physical** **track** **number** **within** **elude**. * , * , 0.

sector n. a region within " a hat t, ac. ** ' =yl. nd / , yli **,
easy to carry out this u-a, silKKn practically . *

g Most of the disks contain s on s a . " * »
me mapping hides this s on e Mty
scotors from elsewhere on tw \ . 'ti.u y "• bu

0 The number of sectors per track is not a constant.

Syllabus Topiu ; ui,)

5.3.3 Disk Attachment

On small systems, disk storage is accessed via VO ports called as host. In a hybrid file system this access is via network storage. In remoted storage, the host is called as host.

5.3.3(A) Host-Attached Storage

Q. Explain host-attached storage.

Local *VO* ports are used to access ~~h~~^h attached storage. Different technology is used by

Desktop PC makes use of an VQ bus **architecture** referred as IDE or ATA. Maximum of two drives per I/O bus is supported by this architecture. In SATA protocol the cabling is simplified, SCSI and fibre channel (FC) is used by high ended servers and workstations.

- SCSI bus architecture has its physical medium as a ribbon cable that contains 50 or 68 conductors. There can be at the most 16 devices supported by SCSI on the bus. In general, the devices contain one controller card in the host called as SCSI initiator. Devices also includes up to 15 storage devices called as the SCSI targets.

The common SCSI target is SCSI disk. The protocol offers the capacity to address up to 8 logical units in each SCSI target. The logical unit addressing is direct or indirect commands to components or a RAID array or components of a removable media library.

FC is a 'Pettiy serial arra... J age Mana m
optical fiber has two Or < or haVi... 8 a erent 24 hil expe... to dominate in the storage-are a 'Pettiy serial arra... J age Mana m
optical fiber has two Or < or haVi... 8 a erent 24 hil expe... to dominate in the storage-are
ou o ra_hal On< inDUCOr U surge switched fabric
funct on oveT
PPC cable. It
address space. This alternative is
U re an I * th@1uis ot
3.3(B) N 8 netl *orki,

5.3.3(B) Network-Attached Storage (NAS)

A network-attached storage device.

itfcj" in UNOJ XX" re C (WO interface
mX 5" «ed by di ' . CIFS TM *!"0">'
the tn' TCP or UDP u uu-j" for 5 n, l ork l" chsd
etworks. Performing RPC over

The network-attached implemented as a R_u unit generally dements the R_u with softw are that PROTOCOL that u/J? interACC NAS is a storage-access

that do? S R P C o Storage access
NAS offers TCP/IP.
tA N to share a *c fOT a n the machinea on a
access similar! storage with simple naming and

having less, rather hand, it is less efficient and
storage option Ormante than TOme direct-attached

ISCS! is the most rCCn! **network**-attached storage protocol \ **l**o
'«iy the Sarpro . “ uses the ff “ lwork protocol lo

5 3 -3(C)Storage-Area Network (SAN)

One limitation of NAS systems is that the storage I/O operations require more bandwidth on the data network. Therefore, latency of network communication increases. This problem can be sensitive in case of large client-server installations where the communication between servers and clients needs more bandwidth along with bandwidth needed for communication among servers and storage devices. Hence, there is competition for bandwidth during communication which can lead to decrease in performance.

- Private network is the example of SAN. It uses storage protocol to connect server and storage unit instead of networking protocol. SAN is very much flexible. Several systems and storage arrays can be attached to SAN and dynamic allocation of storage can also be done to these machines.

A SAN switch permits or restricts access between the Ls and the storage. Consider the example. If host ~~...mn~~ have low disk space then SAN can be configured to offer ffloc StOTa8elo hoSt

Operating System (MU - Sem 4 - IT)

SAN's, it is possible for clusters of servers to share the same storage and for storage arrays to comprise many direct host connections. SANs usually have several more, and less costly ports, than storage arrays.

SyftobU* Topic : Disk Scheduling

5.3.4 Disk Scheduling → (Nov. 18)

Explain the technology of disk scheduling. MU - Dec 2015, 10 Marks

Firstly, we have to consider, how long it takes to read or write a disk block. The time required is determined by considering three factors:

1. **Seek time**: It is the time to move to the proper

2. **ZX-I delay**: The time for the proper sector to rotate under the head.

3. **Actual data transfer time**.

For most of the disks, the seek time is greater than rotational delay and actual data transfer time.

Hence by minimizing mean seek time, it is possible to enhance system performance significantly.

If the processing of requests by disk driver is in First-Come, First-Served (FCFS) basis then it is not possible to optimize seek time.

Yet, another approach is possible in case of heavily loaded disk. When the arm is seeking for one request, other disk requests may be produced by other processes.

Several disk drivers keep a table which is indexed by cylinder number. In this table, all the pending requests for each cylinder chained together in a linked list headed by the table entries.

5.3.4(A) FCFS Scheduling Algorithm

→ (Dec. 16)

Q. Compare the following Disk scheduling algorithms using appropriate example-SSTF, FCFS, SCAN, C-SCAN, LOOK. MU - Dec 2016, 10 Marks

- The first-come, first-served (FCFS) algorithm is the simplest scheduling algorithm.
- This algorithm is basically fair. It general does not offer the fastest service.

Storage Management
R queue with requests for I/O to blocks

Consider

o yim 16, 120, 55, 57

96, is - 35, 122, 16, 120, 55, 57
nuek e = 96, 185, 35, 122, 16, 120, 55, 57

Head start at 51

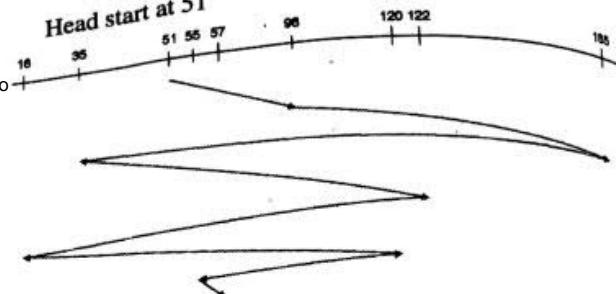


Fig. 53.2 : FCFS disk Scheduling

Initially, the disk head is at cylinder 51. It will then first move from 51 to 55, then to 57, then to 96, then to 120, and finally to 122. The total head movement is 648 cylinders. This schedule is diagrammed in Fig. 53.2.

The larger swing from 122 to 16 and then back to 51 demonstrates the problem with this schedule. If the requests for cylinders 35 and 16 could be serviced together, before or after the requests at 122 and 120, the total head movement could be decreased considerably, and performance could be thereby enhanced.

5.3.4(B) Shortest-Seek-Time-First (SSTF) Scheduling Algorithm

→ (Dec. 16)

Q. Compare the following Disk scheduling algorithms using appropriate example-SSTF, FCFS, SCAN, C-SCAN, LOOK. MU - Dec 2016, 10 Marks

This algorithm services all the requests near to the current head position before moving the head far away to service other requests. The SSTF algorithm chooses the request with the minimum seek time from the current head position. As seek time increases with the number of cylinders pass through by the head, SSTF chooses the awaiting request closest to the current head position.

In our example, the closest request to the original head position (51) is at cylinder 55. Once we are at cylinder 55, the next closest request is at cylinder 57. From there, the request at cylinder 35 is nearer than the one at 16.

96, 50, 35 is served next. Continuing, we service the request at cylinder 16, then 96, 120, 122, and finally 185.

This scheduling method results in a head-to-head comparison of FCFS and Round Robin scheduling. The Round Robin algorithm is considered to be a better method for performance.

Just, ite a shortest-job-f™ (SJF) Kh ~~editi~~ ~~on~~ ~~on~~
 may, "d" Ration of wme reques(s) Reme £
 , a. ~~reques~~ may "7" " time + **Suppose** lh, *
 e tw° ueStS " " queue for linden, j 6 and
 and white servicing .he requesl from
~~10~~ ~~request near~~ >6arrives. 16, a new

- St - ' - : TM
serviced, another request close i₀ 16 cool₁₄ arrive, j_n^g
theory a continual stream of requests near one another
ould arrive. causing the request for cylinder j₈s i₀
wait indefinitely-

One, ee 96. 185,35. 122, 16. 120,55,57

Head start at 51

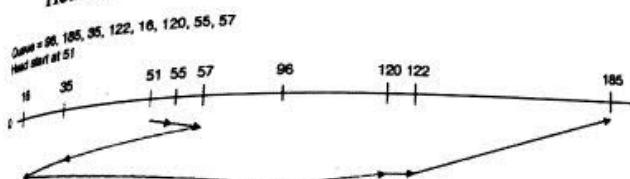


Fig. 5.33 : SSTF disk Scheduling

1.3.4(C) SCAN Scheduling Algorithm

-> (Dec. 16)

7. Compare the following Disk scheduling algorithms using appropriate example-SSTF, FCFS, SCAN, C-SCAN, LOOK (MU - Dec. 2016, 10 Marks).

- The movement of the disk arm in this algorithm is from one end of the disk and goes in the direction of the other end. While doing so, it service requests as arrive at each cylinder, until it arrive at the other end (the disk.

(he disk. The density of requests is very heavy is at the other end of the disk. These requests have also waited for the longer period of time, so why not go there first? That is the consideration of the next algorithm.

fee =96, J85, 35, 122, 16, 120, 55, 57

Head stan at 51

Storage Management

After reaching at the o^{th} end, direction of head move^l $Wh < 8^{cls}$ reversed, and servicing continues. The head ^{re} $Petitively$ scans backward and forward across the disk,

T** arm acts just like an elevator in the building.
Initially it serves all the requests while going up and when reverse the direction towards down, it services the requests on this way. Therefore the SCAN algorithm is also called as the elevator algorithm.

I 1 s consider the same example discussed above. Before proceeding with SCAN to schedule the requests on cylinders 96, 185, 35, 122, 16, 120, 55, and 57, the direction of head movement should be known along with the head's present position (51).

- If movement of the disk arm is towards direction of 0, the head first will service 35 and then 16. After reaching at cylinder 0, the arm will reverse and will go toward the other end of the disk, servicing the requests at 55, 57, 96, 120, 122, and 185,
 - If the incoming request in queue is just in front of the head, then it will be serviced almost without delay; a incoming request just at the back the head need to wait until the arm goes to the end of the disk, reverses direction, and comes back.

Considering an even allotment of requests for cylinders, consider the density of requests when the head arrives at one end and" reverses direction. At this moment, comparatively a small number of requests are right away in front of the head, since these cylinders have just been serviced.

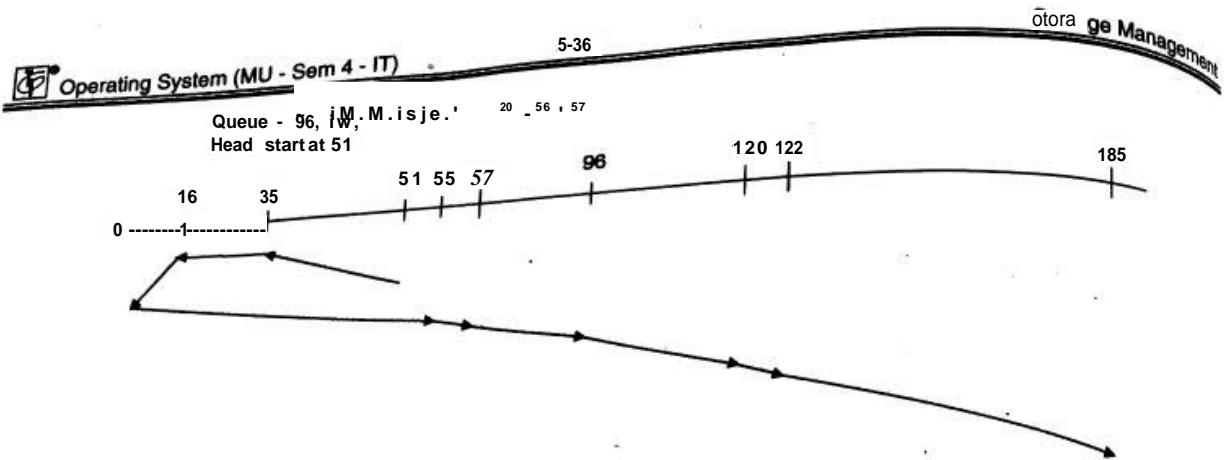


Fig.5.3*4: SCAN disk Scheduling

5.3.4(D) C-SCAN Scheduling Algorithm

Q. Compare the following Disk scheduling algorithms using appropriate example-SSTF, FCFS, SCAN, C-SCAN, LOOK. (Dec. 16)

MU - Dec 2016.7b"

- C-SCAN is called as Circular SCAN scheduling and it is a variation of SCAN. It offers a more uniform wait time. Just like the SCAN, C-SCAN moves the head from one end of the disk to the other, servicing requests along the way.
- When the head arrive at the other end, however, it straight away returns to the beginning of the disk, without any requests on the return trip.
- The C-SCAN scheduling algorithm essentially treats the cylinders as a circular list that wraps around from the final cylinder to the first one.

Queue= 96, 185, 35, 122, 16, 120, 55, 57

Head start at 51

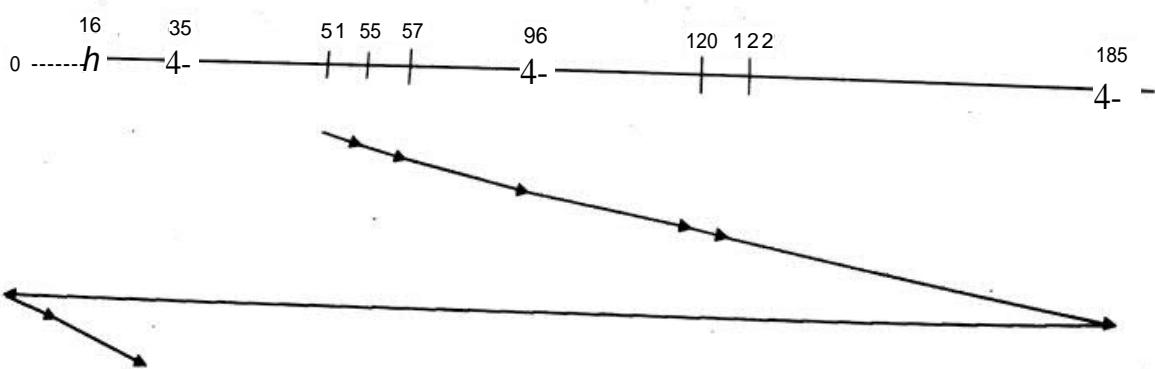


Fig. 5.3.5 : C-SCAN disk Scheduling

5.3.4(E) LOOK Scheduling Algorithm

Q. Compare the following Disk scheduling algorithms

→ (Dec. 16)

US in appropriate example-SSTF, FCFS, SCAN, C-SCAN

MU - Dec. 2016. 10 Marks

Algorithm: SCAN

and v

AN, move the disk across, head width, disk.

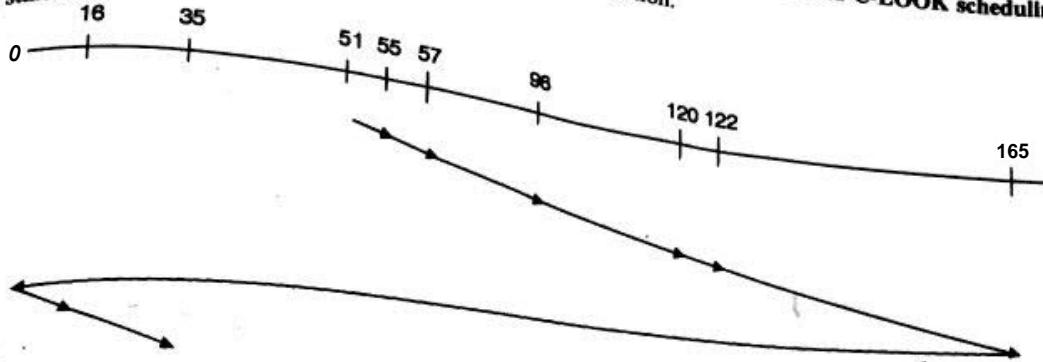
Normally, the arm goes only as far as the final request in each direction. Then, it reverses direction straight away.

5-37

Storage Management

Versions of SCAN and C-SCAN that follow this above pattern are called LOOK and C-LOOK scheduling, because they look for a request before continuing to move in a given direction.

Head start at 51



5.3.6 Examples on Disk Scheduling Algorithms

Fig. 5.3.6 : LOOK scheduling

Example 5.3.1

Suppose that a disk drive has 5000 cylinders, numbered 0 to 4999. The drive is currently serving a request at cylinder 143, and the previous request was at cylinder 125. The queue of pending requests, in FIFO order, is 86, 1470, 913, 1774, 948, 1509, 1022, 1750, 130. Starting from the current head position, what is the total distance (in cylinders) that the disk arm moves to satisfy all the pending requests, for each of the following disk scheduling algorithms.

1. FCFS

2. SCAN

5*rtk)n:

The FCFS schedule is 143, 86, 1470, 913, 1774, 948, 1509, 1022, 1750, 130. The total seek distance is 7081.

The SCAN schedule is 143, 913, 948, 1022, 1470, 1509, 1750, 1774, 4999, 130, 86. The total seek distance is 769

5.3.2

- requests come in to the disk for cylinders 10, 22, 20, 2, 38. A seek takes 6 msec per cylinder move. How much time is for Closest cylinder next algorithm? tyam is at cylinder 20.

Solution:

• Come, First-Served algorithm accepts requests

at time and carries them out in that order,

$$20 - 10 + (22 - 10) + (22 - 20) + (38 - 20)$$

$$+ (40 - 2) + (40 - 10) + (38 - 6)$$

Cylinder Seek « 77' * = 876 ms « Algorithms:

SSF Shortest Seek First B, BORing always handles the closest request time, independently

Total seek time, « 20 - 20 + (22 - 20) + (22 - 20) + (10 - 6) + (6 - 2) + (38 - 2) + (40 - 38) * 6 msec = 6 msec = 360 msec

nearest requests in increasing order handles algorithm handles order until no requests on the way. Afterwards the direction is reversed.

$$\text{Total time: } ((20 - 20) + (22 - 20) + (38 - 22)$$

$$+ (40 - 38) + (40 - 10) + (10 - 6) + (6 - 2)) * 6 msec = 58 * 6 msec = 348 msec$$

Example 5.3.3 MU - Dec. 2014. 10 Marks

Suppose that a disk drive has 5000 cylinders, numbered 0 to 4999. The drive is currently serving a request at cylinder 143, and the previous request was at cylinder 125. The queue of pending requests, in FIFO order, is 86, 1470, 913, 1774, 948, 1509, 1022, 1750, 130. Starting from the current head position, what is the total distance (in cylinders) that the disk arm moves to satisfy all the pending requests, for each of the following disk-scheduling algorithms?

- | | | | |
|---------|----------|---------|---------|
| a. FCFS | & SSTF | C, SCAN | d. LooK |
| OSCAN | I.C-LOOK | | |

- Solution :
- The FCFS schedule is 143, 86, 1470, 913, 1774, 948, 1509, 1022, 1750, 130. The total seek distance is 708.
 - The SSTF schedule is 143, 130, 86, 913, 948, 1022, 1470, 1509, 1750, 1774. The total seek distance is 1745.
 - The SCAN schedule is 143, 913, 948, 1509, 1750, 1774, 4999, 130, 86. The total seek distance is 9769.
 - The LOOK schedule is 143, 913, 948, 1022, 1509, 1750, 1774, 130, 86. The total seek distance is 3319.
 - The C-SCAN schedule is 143, 913, 948, 1022, 14, 1509, 1750, 1774, 4999, 86, 130. The total seek distance is 9813.
 - (Bonus) The C-LOOK schedule is 143, 913, 948, 1022, 1470, 1509, 1750, 1774, 86, 130. The total seek distance is 3363.

Example 5.3.4

None of the disk-scheduling disciplines, except FCFS, is truly fair (starvation may occur).

- Explain why this assertion is true.
- Describe a way to modify algorithms such as SCAN to ensure fairness.
- Explain why fairness is an important goal in a time-sharing system.
- Give three or more examples of circumstances in which it is important that the operating system be unfair in serving I/O requests.

Solution :

- New requests for the track over which the head currently resides can theoretically arrive as quickly as these requests are being serviced.
- All requests older than some predetermined age should be compulsorily placed at the top of the queue, and an associated bit for each could be set to specify that no new request could be moved ahead of these requests. For SSTF, the rest of the queue would have to be reorganized with respect to the last of these "old" requests.
- To avoid unusually long response times
- User requests should have lower priority over paging and swapping. It may be desirable for other kernel terminated I/O, such as the writing of file metadata, to take priority over user I/O. If the system supports real-time process priorities, the real-time processes should be favoured.

Example 5.3.5

Disk requests come into the disk for cylinders 10, 22, 20, 2, 50, 948, 1350, 1500.

40, 6 and 38. A seek takes 6 millisecond per cylinder move. How much seek time is for Closest cylinder next algorithm? Initially arm is at cylinder 20.

Solution ■ First - ComeT and carries them out in that order;

On Total time.
$$((20 - 10) + (20 - 6) + (20 - 2)) * 6 \text{ msec} = 876 \text{ msec.}$$

Cylinder Seek Algorithms: SSF Shortest Seek First

handles the closest request independently on the request, * >

Total time :
$$((20 - 20) + (38 - 20) + (22 - 38)) * 6 \text{ msec} = 60 * 6 \text{ msec} = 360 \text{ msec.}$$

Cylinder Seek Algorithms : Elevator algorithm handles the closest request in increasing (decreasing) order until no more requests on the way. Afterwards the direction is reversed.

Total time :
$$((20 - 20) + (22 - 20) + (38 - 22) + (40 - 38) + (10 - 6) + (6 - 2)) * 6 \text{ msec} = 58 * 6 \text{ msec} = 348 \text{ msec.}$$

Example 5.3.6

Suppose that a disk drive has 5000 cylinders, numbered 1 to 4999. The drive is currently serving a request at cylinder 155 and the previous request was at cylinder 170. The queue of pending requests is 86, 1350, 948, 130, 1500, 50. Starting from the current head position, what is the total distance (in cylinders) that the disk arm moves to satisfy all the pending requests, for each of the following disk scheduling Algorithms.

- FCFS
- SSTF
- SCAN
- C-SCAN
- LOOK
- C-LOOK

Solution :

- FCFS : The FCFS schedule is 155, 86, 1350, 130, 1500, 50

$$= 155 - 86 = 69, 1350 - 86$$

$$= 1264, 1350 - 948 = 402,$$

$$948 - 130 = 818, 1500 - 130 = 1370, 1500 - 50$$

$$= 1450$$

Total head movements are :

$$69 + 1264 + 402 + 818 + 1370 + 1450$$

$$= 5373 \text{ Cylinders}$$

- SSTF : The SSTF schedule is 155, 130, 50, 948, 1350, 1500

Operating System (MUjjem)

$$(4999-1777) = 3222, (4999-130) = 4669,$$

$$(130-80) = 50$$

$$\text{Total head movements are: } 770 + 35 + 74 + 448 + 280 + 27 + 3222 + 4869 + 50$$

= 9775 Cylinders

5. C-SCAN : The C-SCAN schedule is 4999, 0, 80, 143, 913, 948, 1022, 1470, 1750, 1777, J30.

$$\text{It gives } (913-143) = 770, (948-913) = 35,$$

$$(1022-948) = 280,$$

$$(1470 - 1022) = 2, (1750 - 1470) = 280,$$

$$(1777 - 1750) = 27, (4999 - 1777) = 3222,$$

$$(4999 - 0) = 4999,$$

$$(80-0) = 80, (130-80) = 50$$

Total head movements are:

$$770 + 35 + 74 + 448 + 280 + 27 + 3222 + 4999 + 80$$

+ 50 = 9985 Cylinders

Example 5.3.8

Suppose that a disk drive has 200 cylinders, number 199. The initial head position is at 100th track. The queue of pending requests in FIFO is 55, 58, 39, 18, 90, 160, 150, 184. Calculate average seek time for each of the following algorithms,

1. FCFS
2. SSTF
3. SCAN
4. LOOK
5. C-SCAN
6. C-LOOK

Solution :

1. FCFS : The FCFS schedule is

100, 55, 58, 39, 18, 90, 160, 150, 38, 184.

$$\text{It gives } (100-55) = 45, (58-55) = 3,$$

$$(58-39) = 19, (39-18) = 21,$$

$$(90-18) = 72,$$

$$(160-150) = 10, (150-38) = 112,$$

$$(184-38) = 146$$

Total head movements are:

$$45 + 3 + 19 + 21 + 72 + 10 + 112 + 146$$

= 428 Cylinders

2. SSTF : The SSTF schedule is

100, 90, 58, 55, 39, 38, 18, 150, 160, 184

$$\text{It gives } (100-90) = 10, (90-58) = 32,$$

$$(58-55) = 3, (55-39) = 16,$$

$$(39-38) = 1, (38-18) = 20,$$

$$150 - 18 = 132, (160-150) = 10,$$

$$(184-160) * = 24$$

Total movements are:

$$16 + 1 + 20 + 132 + 10 + 24$$

= 248 Cylinders

The SCAN schedule is

SCAN ! - 10, 39, 18, 0, iso, 160, iM

$$100, 5 / 100, 90) = 10, (90 - 58) = 32,$$

$$(58 - 55) = 3, (55-39) = 16,$$

$$(39 - 38) = 1, (38-18) = 20,$$

$$(18 - 0) = 18, (150-0) = 150,$$

$$(160-130) = 10, (184-160) = 24$$

Total head movements are:

$$32 + 3 + 16 + 1 + 20 + 18 + 150 + 10 + 24$$

= 284 Cylinders

4. LOOK - The LOOK schedule is

$$100, 90, 58, 55, 39, 38, 18, 150, 160, 184$$

$$\text{It gives } (100 - 90) = 10,$$

$$(90 - 58) = 32,$$

$$(58-55) = 3, (55-39) = 16,$$

$$(39 - 38) = 1, (38-18) = 20,$$

$$(150 - 18) = 132, (160-150) = 10,$$

$$(184-160) = 24$$

Total head movements are:

$$10 + 32 + 3 + 16 + 1 + 20 + 132 + 10 + 24$$

= 248 Cylinders

5. C-SCAN : The C-SCAN schedule is

$$100, 150, 160, 184, 199, 0, 18, 38, 39, 55, 58, 90$$

$$\text{It gives } (150-100) = 50, (160-150) = 10,$$

$$(184-160) = 24,$$

$$(199-184) = 15, (199-0) = 199,$$

$$(18-0) = 18, (38-18) = 20,$$

$$(39-38) = 1, (55-39) = 16,$$

$$(58-55) = 3, (90-58) = 32$$

Total head movements are:

$$50 + 10 + 24 + 15 + 199 + 18 + 20 + 1 + 16 + 3 + 32$$

= 388 Cylinders

6. C-LOOK : The C-LOOK schedule is

$$100, 150, 160, 184, 18, 38, 39, 55, 58, 90$$

$$\text{It gives } (150-100) = 50,$$

$$(160-150) = 10,$$

$$(184-160) = 24,$$

Total movements are
 $9+5+20+100+10+30 = 230$

(ii) C-SCAN .., 0.37. 54. 43. 40. 20. 0.120.

C-SCAN schedule is 90

150, 180

It gives $(90-80) + (80-57) + (57-54) + (54-45)$
 $+ (45-40) + (40-20) + (20-0) + (120-0) + (150-120)$
 $+ (180-150)$

Total head movements are $10+23+3+9+5+20+20+120+30 = 270$

Example 5.3.12 MU-Dec-201b-10M*±s
 On a disk with 1000 cylinders numbered 0 to 999 compute
 the number of tracks. The disk arm must move to satisfy all
 requests in the queue. Assume the head is moving toward
 the service arm. The queue in FFO order contains requests for the
 head to move to 692, 475, 105, 376.

Perform the computation for the following scheduling algorithms.

j) FIFO S) SSTF \rightarrow SCAN

Solution :

i) **FIFO** : The FIFO schedule is 345,123,874,692,475,105 and 376.

$$\text{Total head movement} = (345-123) + (874-123) + (874-692) + (692-475) + (475-105) + (376-105) = 2013$$

if) **SSTF** : The SSTF schedule is 345,376,475, 692,874, 123 and 105

$$\text{Total head movement} = (376-345) + (475-376) + (692-475) + (874-692) + (874-123) + (123-105) = 1298.$$

(iii) **SCAN** : 345,123, 105, 0, 376, 475, 692 and 874

$$\text{Total head movement} = (345-123) + (123-105) + (105-0) + (0-376) + (376-475) + (475-692) + (692-874) = 1219.$$

Syllabus Topic : Disk Management

5.3.6 Disk Management

The disk management activities of the operating system include disk initialization, booting from disk, and bad-block recovery.

5.16(A) Disk Formatting

If the magnetic disk is new then it is a P' a of a magnenc recording raaterial . Directly on such surface we cannot store data. The surface needs to be divided into sectors so that disk controller is able to perform read and write operations.

The procedure of doing this is called low level or physical formatting. This formatting fills each sector with data structures :

- o Header
- o A data area (usually 512 bytes in size)
- o Trailer

The header and trailer hold the information used by disk controller such as error correcting codes ECC and sector number. ECC is updated with a value calculated from all the bytes in the data. This is done during the controller writes a sector of data at the time of normal I/O, die area.

The ECC is recalculated at the time sector is read and compared with the stored value. If the stored and calculated numbers matches data is correct and corrupted. If mismatch occurs, it indicates that the area of the sector has become corrupted that the disk sector may be bad.

The ECC contains sufficient information, if data have been corrupted, it facilitates the controller to recognize which bits have changed and calculate their correct values should be. It then reports a recoverable soft error. The controller automatically does the ECC processing whenever a sector is read or written.

During the manufacturing process, low level formmation of hard disk is carried out. It helps to manufacture initialize the mapping from logical block numbered defect-free sectors on the disk.

For low level formatting, the number of bytes of data space to leave between header and trailer of all sectors at the time when perform low level formatting. For example, the sizes can be 256,512, and 1,024 bytes. If disk is formatted with larger sector size, fewer sectors can fit on each track.

As a result fewer headers and trailers are written on each track and more space is obtainable for user data.

Some operating systems can handle a sector size of 512 bytes. Operating system keeps its own data structures on disk before it uses disk to store the files. It performs this with following two steps :

- o It partitions the disk into one or more groups of cylinders. Each partition is treated by OS as a separate disk.
 - o Logical formatting. That means creation of file system.
- In order to increase the efficiency, file system groups blocks in chunks called as clusters.

fflz operaung systems give special programs to use a disk partition a, a lar ge sequential array of logical structures. This is somet imes called m da > struc tures. This is called as raw I/O.

Boot Block

S³ bootstrap program is required f_w a ter * . The boot g after it is powered up or it es 4,1 coni Po nenUi Of the system inu s w device controller* and the com regis tery and then starts the operating mem ory program, be OS kernel * * * t kerne! into memoiy, and ju mps to " d,sk addre0 <> start * he Pera,in 8- tem execution.

The Read Only Memory (ROM) does not require initialization and is at a fixed location that the processor can begin executing when powered up or down. Therefore bootstrap is stored in ROM.

* use Of read only feature of ROM; it is not corrupted by application » a computer vnts. diff Ufes chan ROM hardware chips.

Lafore. most s > 8,6,ns store a 'Mil bootstrap loader program in the boot ROM. Wh,Ch inVokes a,ld bri "g full Ldstrap program from disk into main memory. The Ldified version of full bootstrap program can be simpJy written onto the disk.

The fixed storage location of full bootstrap program is in the "boot blocks". A disk that has a boot partition is called a boot disk or system disk. Boot ROM code gives instruction to controller for reading the boot blocks into memory (no device drivers are loaded at this point) and then starts executing that code.

The full bootstrap program is more sophisticated than the bootstrap loader in the boot ROM; it is able to load the complete operating system from a non-fixed location on disk and to start the operating system execution. Yet, the full bootstrap code may be small.

5.3.6(C) Bad Blocks

Failure of the disk can be :

- Complete, means there is no way other than replacing the disk. Back up of content must be taken on new disk.
- One or more sectors become faulty.
- After manufacturing, the bad blocks exist. Depending on the disk and controller in use, these blocks are handled in a different ways.

If the disks are with IDE controllers then these are simple disks. In these disks bad blocks are handled "smoothly". For example, the MS-DOS format command carry out logical formatting and as a part of the process

Storage Many mgl

scans the disk to discover bad blocks. If format discovers a bad block, it writes a special value into the related I A entry. These blocks go bad during normal operation, a program (such as chkdsk) must be run to look for these blocks and to lock them away. Data that resided on these blocks generally are lost.

IVn. The 5051 disks are used in high-end PCs and most workstations and servers, are more intelligent about recovery of bad-block. The controller preserves a record of blocks on the disk. The list is initialized during the low-level formatting at the factory and is updated over the life of the disk. Low-level formatting also sets aside spare sectors not visible to the operating system. The controller can be told to replace each bad sector logically with one of the spare sectors. This is called as sector sparing or forwarding.

A usual dealing with bad-sector is like follows :

- The OS attempt to read Logical block 67.
- The controller computes the ECC and concludes that the sector is bad. The same result is reported to OS.
- When the system is rebooted, a special command is run to inform the SCSI controller to replace the bad sector with a spare.
- **Afterward, at whatever time the system requests logical block 67, the request is translated into the replacement sector's address by the controller.**

Syllabus Topic : RAID Structure

5.3.7 RAID Structure

5.3.7(A) RAID Levels

Q. Explain various RAID levels.

- The performance of CPU has been rising immensely over the past decade, approximately it is getting increasing by 100 percent after every 1.5 years. On the contrary disk performance is not increasing with this speed. Now days average seek time of the disk is near about 10 msec, which is very less compare to average seek time of 50 to 100 msec on minicomputers in the decade of 1970.
- If in any other industry, performance is enhanced by 5 to 10 percent in two decades, it would be considered as great achievement. However, in the computer industry it is an discomfiture. As a result, the difference between CPU performance and disk performance has grown to be much larger in due course of time.
- In order to get faster processor performance, parallel processing is being used to a greater extent. Over the years people thought that parallel I/O might be an excellent scheme as well. In research paper of 1988, Patterson et al. recommended six specific disk

Opting System (MU. Sggj)

organizations that could be utilized to gain better disk performance, reliability, or both.

This proposal was soon agreed by most organizations and has made possible a new class of RAID. Patterson et al. had given "Inexpensive Disks, but RAID as Redundant Array of Inexpensive Disks" as industry soon changed the definition of "Independent" instead of "Inexpensive".

The fundamental thought behind a RAID is to have a bunch of disks beside the computer, use a RAID controller card, copy the data on RAID and then carry on usual operation. The intention was a KIPS system to appear a single large expensive disk to the operating system but have improved performance and reliability.

As SCSI disks are characterized by better performance, low cost, and the capability to include up to 7 hit a single controller/ 15 for wide SCSI), it is usually most RAID contain a RAID SCSI controller in addition to a bunch of SCSI disks which operating system can think as a single large disk. In this fashion, no software modifications are necessary to make use of the RAID, a great selling point for a majority of system administrators.

Besides coming out similar to a single disk to the software, every RAID have the feature of distributing the data over the drives, to permit parallel operation. Many varieties of methods for achieving this were defined by Patterson et al., and they are at present recognized as RAID level 0 through RAID level 5. There are only six distinct organizations possible.

RAID level 0 is demonstrated in Fig. 5.3.7(a). It gives the impression of the virtual single disk simulated by the RAID as being split up into strips of k sectors each, by means of sectors 0 to $k-1$ being strip 0, sectors k to $2k-1$ as strip 1, and so on. Each strip appears to be a sector for value $k=1$; a strip is two sectors if $k=2$, etc. The RAID level 0 collection writes successive strips over the drives in round-robin manner, as shown in Fig. 5.3.7(a) for a RAID having four disk drives.

When the data is distributed like this on more than one drives, the operation is called as striping. If data block comprises four consecutive strips beginning at strip boundary and software gives a command to read this data block, the RAID controller will split this command up into four separate commands, one for each of the four disks, and have them run in parallel. Thus parallel I/O is achieved and software remains unaware about it.

If the large numbers of request are present then DMR

level 0 works better. If number of drives times the strip size is less than the request, some of the drives will

receive more than one requests, with the intention that when they comply with the first request they initiate the second.

The controller decides to divide the requests up and provide the right order and then collect the appropriate disks in right manner. Performance is simple. RAID 0 and the only with operating systems that support it is DOS. It suits a user to demand K but there is no increase in performance due to parallelism.

Another drawback of this organization is that the probability of failure free operation is potentially poor than having a single large expensive disk (SLED). If a RAID includes four disks, having mean time to failure (MTTF) of 20,000 hours for every disk, about once every 5000 hours a drive will crash and all the data vanished. A single drive with a MTTF of 20,000 hours is 100 times more reliable. Since redundancy is absent in design, it is not actually a true RAID.

The next alternative, RAID level 1, demonstrated in Fig. 5.3.7(b), is an exact RAID. In this, each disk is present, so total 8 disks are present. Of these 8 disks, four are primary disks and four are mirror disks. The strip is written twice for each operation. A read can be carried out on any of the drives by assigning and distributing the load over more drives.

As a result, write performance is comparable to a single drive, but read performance can be improved neatly by double. Fault tolerance is outstanding: if a drive fails, the copy is just used in place. Recovery includes simply fixing a new drive and copying the complete backup drive to it.

The RAID levels 0 and 1 operate with strips of sectors. Whereas RAID level 2 supports the working on either word basis or even on a byte basis. Consider division of each byte of the single virtual disk into a two 4-bit nibbles, then inserting a Hamming code to each one to build a 7-bit word having bits 1, 2 and 4 parity bits.

Consider again Fig. 5.3.7(c) in which the seven drives are synchronized with respect to arm position and rotational position. Then it would be possible to write one bit per drive out of the 7-bit Hamming coded word.

This idea used by the Thinking Machines' CM-1 computer. It considers 32-bit data words and append 6 parity bits to form a 38-bit Hamming word, in addition to one additional bit for word parity, and distribute each word over 39 disk drives.

The total throughput was enormous, because it was possible to write 32 sectors worth of data in one second writing time. The crash of one drive did not lead to problems, as loss of a drive cause the loss of 1 Mbit in each 39-bit word, which Hamming code can manage easily.

I have raw knowledge of the above, and the main point is that, in this, the drives need to be rotated sequentially.

of drives (regardless of the controller, the overhead is just a few percent, and every bit is reliable at the HAJD level 3 level,™). A simpler version of RAID 1 is

/x: i-sci- iT a>«
/eve, 2. «« -nve. shoula Party da(a P-ee
because mdmduiU data Like 2." U

thus are synchronized and distributed over

Stripe

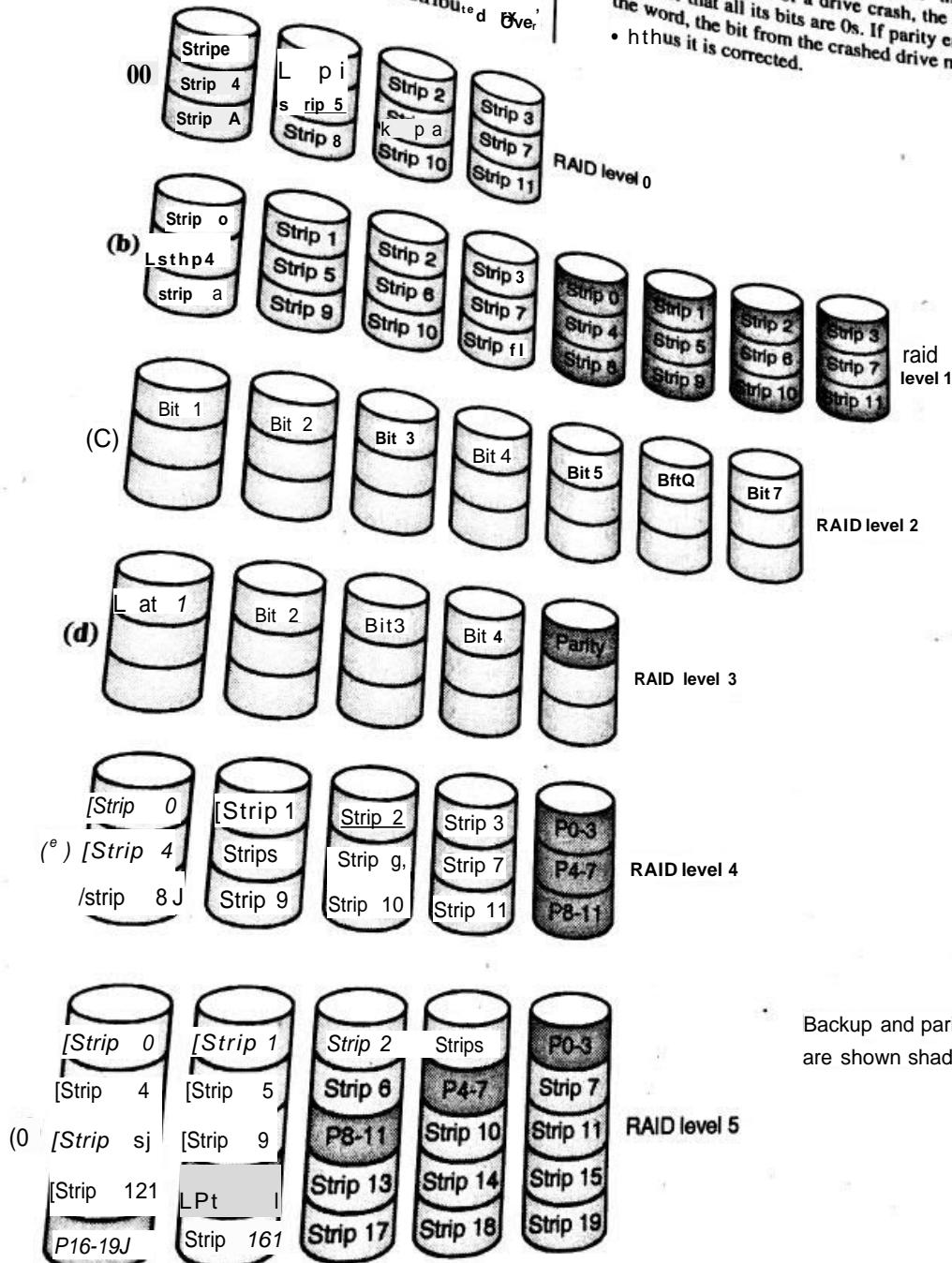


Fig. 53- 7: RAO ,evels 0 thr0Bgh 5

rates, the total separate I/O requests per second

Despite the trend to use MLC RAID levels, 2 to 3 P_{MLC}^{*} are the most common.

they can handle is no better than for a single drive.

Operating System (MU-Sorru)

- In RAID levels 4 and 5 do not need synchronization. In RAID level 4 consider individual words similar to RAID level demonstrated in Fig. 5.3.7(e). Parity is written onto an extent 0, and here a strip-for-strip is carried out. The EXCLUSIVE OR operation is carried out in strips together. This operation results in length k bytes.
- In case of failure of drive, the parity drive is recovered by recomputing from the parity drive but design defends against the failure. If a sector is altered, it is essential to make a write to that be rewritten.
- On the other hand, it may read the old user data and the old parity data and compute again the new parity from them. Still with this optimization, a small update needs to have two reads and two writes.
- As a result of the huge load on the parity drive, it may turn into a bottleneck. This bottleneck is removed in RAID level 5 by splitting the parity bits evenly over all the drives, in round robin manner, as illustrated in Fig. 5.3.7(f). Yet, in case of a failure of drive, rebuilding the contents of the failed drive is a complex process.

Syllabus Topic : Stable-Storage Implementation

5.3.8 Stable-Storage Implementation

- For the write-ahead log, stable storage is required. The information stored in stable storage remains permanent and never lost. In order to implement this stable storage, it is necessary to replicate the information over multiple disks with independent failure modes.
 - The update writing should be carried out in a manner that promises that a failure during an update will not damage all the copies. It should also guarantee that during recovering from a failure, all copies must remain consistent and with correct value, though another failure takes place during the recovery.
- One of the following can happen due to disk writes.

- Successful completion.** The writing of data on disk correctly done.

- Partially written.** If a failure occurs during the write operation, hence, some of the sectors may have been corrupted. At the time of the failure occurred prior to the write, so the previous data value, which is write-32ns unchanged.

- If during a write operation, a failure occurs, then the system must keep two physical blocks for each logical block. Following are the execution steps of each logical block:

1. **Initial write.** The first write completes successfully, the same physical block.
2. **Second write.** The second write completes successfully, and the operation is completed.

- As part of recovery procedure from failure, of each pair of physical blocks is carried out. If the two blocks are identical, then no further action is needed. If one block has a detectable error, then its contents are replaced with the value of the other block. If detectable error does not exist in both blocks, but both differ in content, then it is necessary to replace the content of the first block with that of the second.

- This recovery procedure gives the guarantee that, a write to stable storage either succeeds completely or results in no change. This procedure can be extended with no trouble to permit the use of an arbitrarily large number of copies of each block of stable storage.

Syllabus Topic : Tertiary-Storage Structure

5.3.9 Tertiary-Storage Structure

5.3.9(A) Tertiary-Storage Devices

- Q. Explain various Tertiary-Storage Devices.

The main characteristic of tertiary storage devices is its low cost. Following are the examples of tertiary storage devices.

Removable Disks

- Floppy disk is the example of removable magnetic disk. The storage capacity of floppy disk is only

device. Typically, the tape drive then is reserved for the exclusive use of that application until the application exits or closes the tape device.

- As tape drive is presented as raw device, OS does not offer file system services. The application must make a decision about how to use the array of blocks. The application decides its own rules for organizing a tape, a tape full of data can usually be used by only the program that created it. For example, although a **backup tape** holds a list of file names and their sizes followed by the file data in that order, it is not easy to use the tape.
- **Disk drive** performs operation like `read()`, `write()` and `seek()`. Instead of `seek()`, tape drive uses `locate()` operation. The `locate()` operation of tape is more accurate than `seek()` operation as tape gets positioned to correct logical block instead of entire track. Many number of tape drives supports `read-position()` operation that returns the logical block number where the tape head currently exists. Many tape drives also support `space()` operation for relative motion.

File Naming

- On personal computer, file name contains name of the drive followed by path name. In case of removable disk, it does not mean that if drive containing the cartridge at some time in the past is known means how to find the file is known. If every removable cartridge had a unique serial number, the name of a file on a removable device could be prefixed with the serial number, but to make sure that no two serial numbers are the same would require each one to be about 12 digits in length. It would then difficult to remember 12 digits file number as a file name.
- In case of writing the data on a removable cartridge on one machine and then use the cartridge in another machine, the problem will become more complex. If both machines are same and with same kind of removable drive, then only trouble is to find the contents and data layout on the cartridge. If both machines are not of same types or drives are different, many other problems can come up.

Machines with different architectures have their data representation in different formats. Usually in modern operating systems, this name-space problem is

addressed for removable media. TX

on applications and users to discover how to accept data. But, only some kinds of

data are so well standardized that all machines use them in similar way.

Hierarchical Storage Management

A hierarchical storage management system facilitates for changing the cartridge in a tape or disk drive without intervention. A hierarchical storage system extends the storage hierarchy beyond main memory and secondary storage to include tertiary storage. Tertiary storage is implemented as a jukebox of tapes with removable disks. This storage hierarchy is larger, cheaper, and slower.

The normal method to include tertiary storage is to widen the file system. Small size and frequently used files resides on disk, whereas large and old files having no active use is stored to the jukebox. In some file-archiving systems, the directory entry for the file keeps on existing, but the contents of the file no longer holds space in secondary storage.

Now day, hierarchical storage management is set up in installations with large volumes of data that are used rarely, at irregular intervals, or sometimes. Existing work in HSM comprises extending it to offer full **information life-cycle management (ILM)**. In this case, data go from disk to tape and back to disk, as required, but are deleted on a schedule or in accordance with policy.

5.3.9(C) Performance

Tertiary storage performance is measured in terms of speed, reliability, and cost.

Speed

Bandwidth and latency are two metrics to measure performance of tertiary-storage devices. **Bandwidth** is measured with bytes per second. The bandwidth which is steady is the average data rate in a huge transfer specifically, the number of bytes divided by the transfer time. The effective bandwidth considers the average over the complete I/O time, together with the time for locate() and any

switching time in a jukebox. > an cartridge switching time in

Essentially, the sustained bandwidth is the overall data rate given by the drive.

The sustained bandwidth is usually a bandwidth of a removable disk which are slowest, the bandwidth ranges from a few megabytes per second

and for fastest, it is 40 MB per second. Tapes encompass a same range of bandwidths, from few megabytes per second to over 30 MB per second. Latency is second metrics, 0 measure performance.

which took performance is then better ft, Pa to *»

considerably if we remove it 58 late, cy *Mf
bal ed. then drive must stop 5p I. needed t
ann !S required to swi
and then the drive should spin * ** £

11* up ***>**! tit ** carri
and on-access time within one longer £
pitching disks in a jukebox leads As . re *
w h performance penalty. a con >Paratival y

file robotic-arm time is same for d ,
etching the tape, old tape requires re /** Befo *
ran * e lected . The dn « required fo J ng it
sb out 4 minutes. Several seconds are Opera "TM is
loading of new tape in drive to get ? Uired after
prepare for I/O. In case of performance e o ^{way and}
a jukebox, the bandwidth and latency I t drives

But seem to have a terrible bottleneck oo reasonable

Reliability

Reliability is important metric tn "feasure the performance. Removable magnetic disks " e, o extent is less reliable than are fixed hard disks. In case of removable disks, the cartridge is probably more exposed to damaging environmental conditions; for example, dust, changes in temperature and humidity, and mechanical forces like shock and bending.

Optical disks are very much reliable. In this disk, layer storing the bits is protected by a transparent plastic or glass layer. The reliability of magnetic tape depends on type of drive and it differs largely. Some low-cost drives wear out tapes if used for few number of uses. Whereas, some other type of drives permits millions of reuses. The magnetic-tape drive head is a weak spot compared to magnetic-disk head.

Cost

Storage cost is also important factor that is necessary to consider. Removable media lower the overall storage cost. Though it is somewhat expensive to create a removable cartridge, the cost per gigabyte of removable storage may well be lower than the cost per gigabyte of a hard disk. The reason is, the cost of one d-----

Storage Management

Z' 8 * 1 With tow CM1 ot many removable

** has ; cost ncr m ^{guylte} magnetic
has has more ihul four orders of

No only th ^{ree} orders for primary memory

No way, prim arv 8 Ude.

storage b> "racto OUOTtt? mwt cor, y ora disk

fall * more speedily for disk drives than for tape

magnetic, tSAfove Jre be price p " megabyte of .

the tape drive is 20in o f * lape cartridge without
small- and med . 0 ** n "y saTM • A » result,

S "W cost than diskTM 11 Ub TM les ha m ; a w ^h " landKK *ysteTM With equal capacity.

Swap-Space Management

5.3.10 Swap-Space Management

Q. Explain management of swap Space in detail

the opera bl g e TM agemen to * level job performed by
hold all >h/,, S ma TM memOT y » •>« sufficient to
disk Rnarp J*** 8 simultaneously, virtual memory uses
access k s aCCCSsing the disk s much slower than
reduce system performance. The main aim towards the
Jign, and implementation of swap space is to offer the
best throughput for the virtual memory system.

5.3.10(A) Swap-Space Use

- The use of swap space is carried out in different ways by various operating systems. The way operating system uses swap space merely depends on the type of the memory-management algorithms in use. Such as, Swap space may hold complete process image, with code and data segments.
- Paging system may store the swapped out pages from main memory. The needed size of swap space differs based on the amount of physical memory, the amount of virtual memory it is backing, and the manner in which the virtual memory is used. It can range between few megabytes of disk space to gigabytes.
- It is better to keep sufficient amount of swap space that is required. Otherwise operating system will be forced to abort the processes as system may run out of swap space. There can be chances of system crash due to its running out of swap space.

- Many operating systems and Linux permit the use of multiple swap spaces. These swap spaces are maintained on separate disks. As a result, swap space can be spread over the system's I/O devices.

5.3.10(0) Swap-Space Location

- A swap space can be part of the normal file system, or it can be in a separate disk partition. If it is just a large file in the file system, normal file-system routines can be used to create it, name it, and allocate its space. This scheme is simple to implement but is not efficient. In this approach, navigation of directory structure and the disk-allocation data structures consumes more time and also extra disk accesses are needed. External fragmentation can very much increase swapping times by forcing multiple seeks during reading or writing of a process image.
- Instead, swap space can be created in a separate raw partition. In this approach, a separate swap-space storage manager allocates and reclaims the blocks from the raw partition. The algorithms used by swap-space manager are better in terms of speed. For these algorithms, speed is important instead of efficiency as there is frequent access to the swap space. As data remains for short period in swap space, internal fragmentation is accepted.
- This is not the case with file system. Adding extra swap space needs repartitioning the disk, this requires to move the other file-system, partitions or destroying them and restoring them from backup or adding another swap space in another place. Linux supports both schemes: in raw partition and in file system.

Syllabus Topic 1/0 systems

5.4 I/O systems

5.4.1 Overview

- As we know that I/O is one of the main functions of an operating system and is used to control the entire computer's Input/Output devices.
- Basically it should have the following features, catch interface, commands to the bus. Provide an interface between the devices and the rest of the system that is simple and easy to use.

The computers operate with many kinds of devices. As we know that it includes transmission devices (network cards, modems), and human-interface devices (screen, keyboard, mouse).

a) I/O Devices

Q. Explain different

o, / OdeMicro

Generally when computer system uses VO devices a external devices can be grouped into three categories -

1. Human, "d b k , " * SUI, ab , e Co mnuni "Ung with the computer user.

Examples : PTM terminals, keyboard, and mouse.

2. Machine readable : It is suitable for communicating with electronic equipment.

Examples : Disk drives, USB keys, sensors, controllers, and actuators.

3. Communication : It is suitable for communicating with remote devices.

Examples : Digital line drivers and modems.

- Computers system operates with many kinds of devices. It includes storage devices (disks, tapes) transmission devices (network cards, modems), and human-interface devices (screen, keyboard, mouse).

- A device communicates with a computer system by sending signals over a cable or even through the air. The device communicates with the machine via a connection point termed a port (for example, a serial port). If one or more devices use a common set of wires, the connection is called a bus.

- When device X has a cable that plugs into device Y and device Y has a cable that plugs into device Z, and device Z plugs into a port on the computer. such type of

Zu DUS. is called a daisy chain. " > w

5.4d(B) Differences between I/O Devices

Q.

SZ.

different Parameters distinguish 1/0

Devices can be distinguished with respect to following parameters.

1. Data transfer rate
3. Controlplexiyy
5. Data representation
2. Device application
4. Data transfer unit
6. Errors

5-5
Data transfer rate : Data transfer rate of different devices can be different and can be of several orders

Zsj"»fo Topic : Overview MO

< I/O Hardware

Device Controllers

- * Kput/outiw' U*** S COnsist of a mechanical com...
•» ele " co " ~t. A co ~~ntent~~ r .s a
section of electronics that can operate P o R , a hUS
Dradevice-

A serial-port controller is an example of a device controller. It is a single chip in the computer, controls the signals on the wires of a serial port.

The implementation of SCSI bus controller is done most of the time as a circuit board which separately plugged into the computer.

It normally includes a CPU, microcode, and some private memory to facilitate it to process the SCSI protocol messages. The SCSI bus controller is integral part of the some of the devices.

- Total four registers are present in I/O port. These are as follows :

- o Status register
 - o Control register
 - o Data-in register
 - o Data-out register

The host can read the hits in status register. These bits show states like :

- o The current command is completed or not
 - o A byte exist in data-in register or not which is to be read
 - o If device error present

The host writes the control register to command or to alter the mode of a device. A particular

bit is the control -1
fundamental orange Management

Por ln <<C >> "PlexC* Wrtal Pon * cW << "b,u>

Some of the T are for selection of one of the speeds provided by the serial port.

host gets input by reading data register and writes the data out register in size of data block.

is became semm of ir W and output a f ChipS qan store 8Cv eral bytes to tOre a small biTM / FIFO chip WORK like buffer receive those dau. until the device or host

5 - 4 Wolling

Q. Write note on "Polling".

- Inco
and T c' ProlOC0 , fot totraktion between the host
handshak- can U infiltrate , but basic
Shaking notion is simpk .
The controller indicates its slate through the busy bit in
status register. (Recall that to set a bit means to
*nte a 1 into the bit, and to clear a bit mean to write a
0 into it.)

When controller is busy in operations, it sets a busy bit. This busy bit is cleared when it is not working and ready to accept the next command. The host signals its desires by means of the command-ready bit in the command register.

- This command-ready bit is set by host when a command is available for the controller to execute. ■

For this example, the host writes output through a port, direct to a file, as follows:

1. Until bit becomes clear, the host constantly reads the busy bit.
 2. In the command register, the host sets the write bit and then writes a byte into the data-out register.
 3. Command-ready bit is set by host,
 4. When the controller notices that the command-ready bit is set, it sets the Busy.
 5. After this, the command register is read by control! and sees the write command.

6. In order to get the byte, it reads the data-out register and perform the I/O to the device.
7. The different bits are then cleared as : The controller clears the command-ready bit, clears the error bit in the status register to specify that the device I/O succeeds and clears the busy bit to specify that it is finished.

If the status of host is busy-waiting or polling; It reads the status register repeatedly until the busy bit turns out to be clear. **This means it remains in loop until busy bit clears.** If the controller and device speed is better, this method seems to be practical one. If host has to wait for longer period of time; it should perhaps go for another task.

5.4.2(C) Interrupt Handler

Q. What is interrupt? Explain the tasks carried out by interrupt handler.

- The CPU hardware contains interrupt request line that the CPU senses after executing every instruction. When the CPU detects that a controller has asserted a signal on the interrupt request line, the CPU saves a small amount of state, such as the current value of the instruction pointer, and jumps to the interrupt-handler **routine at a fixed address in memory.**
- The interrupt handler finds out the reason of the **interrupt and carries out the required processing.** After this the interrupt handler executes a return from interrupt instruction to return the CPU to the execution state before the interrupt.
- The device controller first asserts a signal on the interrupt request line. It does so to raise the interrupt. **The CPU catches this raised interrupt and sends out to the interrupt handler** by giving the interrupt handling

="X,TcT,"
asynchronous event, for example a device controller becoming ready for service.

- The modern operating system should support for following interrupt-handling features.

The priorities to process

- It should have efficient

We private interrupt handler to send out to the

out going adjustment, device driver which

A multilevel interrupt should be supported. It helps the operating system to differentiate between high and low priority interrupts so that it can respond with the appropriate degree of necessity.

modern computer hardware, CPU and the interrupt controller hardware offers alt above three features. CPUs have two interrupt request lines. One is still maskable interrupt, which is reserved for events such as unrecoverable memory errors. The second interrupt, is also maskable.

It can now write the CPU before the execution of critical instruction sequences that must not be interrupted. The maskable interrupt is used by device controllers to request service. This address is an offset in a table called the interrupt vector. This vector contains the memory addresses of specialized interrupt handlers.

The purpose of a vectored interrupt mechanism is to reduce the need for a single interrupt handler to service all possible sources of interrupts to determine which one needs service.

The interrupt mechanism also implements a joining of interrupt priority levels. This mechanism enables the CPU to defer the handling of low-priority interrupt without masking off all interrupts, and make it possible for a high-priority interrupt to preempt the execution of a low-priority interrupt.

The interrupt mechanism is also used to handle a wide variety of exceptions, such as dividing by zero, accessing a protected or nonexistent memory address, attempting to execute a privileged instruction from user mode.

A system call is a call that is called by an application to invoke a system service. The system call checks the given by the application and then executes arguments to kernel, special

interrupt, or 'X' inStLIC,IO

to be used with the manage 0.6 flow of efficiently, we need to disks are to be used. Previous one completes the next soon the

Consequently, the kernel is implemented by code that completes a disk

The i/o status is a pair of interrupt handlers.

then cleared by high-priority handler. It begins the subsequently

user level ^{1/0} implemented .. "Pl *
by Co Pyin8 da,a fr om ko he com
han han
lives in space. It then invoke, J" Mera'S
application on the ready queue he *hed ull.r " he
fundamental interrupt meek .

This CPU to respond to an asynchronous event, as when a service Permits the
device controller become, ready for service
0 * modem operating system, however
0 * indicated interrupt handling
still, specified interrupt handling
hi. yet
The ability to postpone internal processing
0 critical processing

0 An efficient technique to dispatch to the proper interrupt handler for a device with all the devices to see which first raise the interrupt- Ming d the

0 The multilevel interrupts, so that the operating system can differentiate between high priority interrupts and can respond with the appropriate degree of ~~urgency~~ ^{urgency}.

54 p) Interrupt Service Routine (ISR)

Interrupt means event, which invites attention of the processor on occurrence of some action at hardware or software interrupt instruction event. In response to the interrupts the routine or program, which is running presently interrupts and an interrupt service routine (**ISR**) executes. **ISR is also called device driver in case of the devices and called exception or signal or trap handler in case of software interrupt.**

- Processor executes the program, called interrupt service routine or signal handler or trap handler or exception handler or device driver, related to input or output from the port or device or related to a device function on an interrupt and does not wait and look for the input ready or output completion or device-status ready or set.

- The programmer may register Interrupt Service Routines (ISRs) to handle hardware interrupts, routines are not independent threads, but more like ~~signals~~ Any currently executing thread is suspended by an interrupt, and the ISR called. The ISR runs on a hardware interrupt stack only if the hardware interrupt is pending. Otherwise, the ISR stack frames are pushed onto the stack of the interrupted thread.

example, in implementing a time-slice scheduler. However the ISR, although more flexible than in typical OSs, still would not be as useful as the tool that allows us to do this. Um, the synchronization especially for CTbSync is provided by ISR routines.)

5.4 "bh, J ISR routines.)
Q: — **ory** Accent (DMA)
Explain, etc.

When an external device requests data transfer, a block of data directly between an external device and the primary memory, without involving the CPU. Rolling or interrupt software makes use of DMA.

- **Disk** involves the transfer of large size information in a **singl.** operation. This transfer of information **be**tween disk and main memory involves DMA operation. When used together with an interrupt, the processor is conveyed only after the whole block of data has been moved.

Information is transferred in the form of bytes or words and for these it must supply the address of memory location and every bus signals that control the data transfer. Communication with a device controller is handled through a device driver. Device drivers are part of the operating system.

- The operating system provides a simplified view of the device to user applications (e.g., character devices vs. block devices in UNIX). In some operating systems (e.g., Linux), devices are also accessible through the /dev file system.

In some cases, the operating system buffers data that are transferred between a device and a user space (X S network buffer). This usually improves performance but not always. The DMA controller has access to the system's memory, but it is not independent of the CPU. In Fig. 3.4.1,

the CPU " 8th processor can -d or write. These of registers dw Kjistelj a byte renters compn^x a memory address several control registers countregister, ands gi

- The function of the control registers is to detail about the facts like I/O port to utilize, reading from or writing to the I/O device, the unit of data or information transfer (in terms of byte or word at a time), and the total bytes to transport in one burst.

Without DMA, let us see how disk read occurs. Initially the controller reads the block from the drive serially, bit by bit, until the complete block is in the controller's internal buffer.

Next, it calculates the checksum to confirm that no read errors have occurred. Then the controller causes an interrupt.

- After this, operating system can read the disk block from the controller's buffer a byte or a word at a time by executing a loop.
- In one byte or word is read from controller device register and storing it in main memory.
- With DMA, the process is different. First the CPU programs the DMA controller by setting its registers so it becomes aware about what to transfer where (step 1 in Fig. 5.4.1). It also instructs to the disk controller to read data from the disk into its internal buffer and confirm the checksum. DMA can start only when valid data are in the disk controller's buffer.

- The DMA controller begins the transfer by issuing a request over the bus to the disk controller (step 2). The disk controller remains unknown and does not care

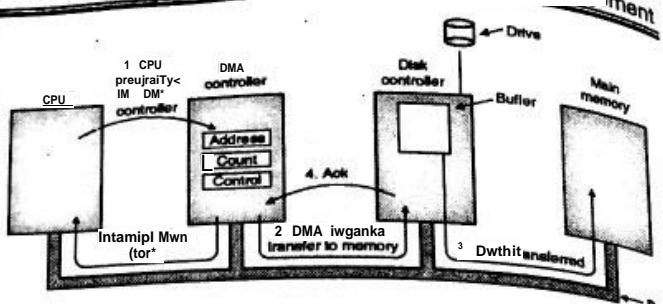
whether read request came from the CPU or from a DMA controller.

Naturally, the memory address to write to is on the bus' address lines so when the disk controller fetches the word from internal buffer it knows where to write it. The write to memory is another standard bus cycle (step 3).

When the write is ends, the disk controller sends an acknowledgement signal to the disk controller, also over the bus (step 4). The DMA controller then increments the memory address to use and decrements the byte count.

If the byte count is still greater than 0, it repeats until the count reaches 0.

At that time, the DMA controller interrupts the CPU to let it know that the transfer is now complete. When the disk block to memory - it has been copied there. This shows DMA operation.



Fig*5.4.1 : DMA transfer operations

Syllabus Topic : Application I/O Interface

5.4*3 Application I/O Interface

- Due to interfaces, operating system treats I/O devices in a standard and uniform way. In VO related part of kernel module, the device-driver layer hides the differences among device controllers from the subsystem of the kernel, much as the I/O layer encapsulate the behavior of devices in a few classes that hide hardware differences across applications.
- If we make the input-output interface of the hardware, then the job of the operating system developer will become simple.
- Every type of operating system has its own interface for the device-driver interface. It is the difference between device-hardware manufacturers.
- A given device may ship with multiple device drivers for instance, drivers for MS-DOS, Windows 95/98, Windows NT/2000, and Solaris. Devices vary on many dimensions.
 - Character-stream or block device** : A character device transfers data one by one.
 - Block device** : A block device transfers a block of bytes as a unit.
 - Sequential device** : A sequential device transfers data in a fixed order determined by the user. Whereas the user of a random-access device can instruct the device to any of the available data storage locations.
 - Synchronous device** : A synchronous device transfers data with predictable response times. An synchronous device exhibits either predictable or unpredictable response times.
 - Asynchronous device** : An asynchronous device exhibits unpredictable response times.
 - Shared device** : A sharable device can be used by several processes or threads.

Speed of operation : Device speeds range from a few bytes per few bytes per second to a few kgab bytes per second. Some devices support only one data transfer but others support both input and output. Some devices support only one data transfer but others support both input and output.

Block or Character Device

4(A) CBB

Block device interface holds information about 8 disk drives and other block devices. The device must recognize the device type if it is a random access device. It can have a seek command to determine the transfer next. These devices are citation through file system interfaces.

Single system management Database management system, LS** the block device as a simple "ear block" application carries out its own buffering then raw device directly to the application.

It helps to step out the OS to carry out some file operations. When cached as raw UO, OS will not be led to perform these operations. UNIX allows a mode of operation on a file that disables buffering and firing called as Direct I/O.

It is advantageous to have memory-mapped file access on top of block-device drivers. A memory-mapped interface offers access to disk storage via an array of bytes in main memory.

The system call that maps a file into memory returns the virtual memory address that contains a copy of the file. The real data transfers are carried out only when required to satisfy access to the memory image.

As the transfers are handled by the same means as that used for demand-paged virtual memory access, memory mapped I/O is efficient. The example of a device that is accessed through a memory-mapped interface is keyboard.

The basic system calls in this interface are used to get() or put() one byte at a time. The interface offers line-at-a-time access with buffering and editing services.

W) Network Devices

Network VO and disk VO differ in terms of performance and addressing.

The majority of operating systems available in many of the operating systems. The interface which is selected when the application is called, to be received by the socket having a packet waiting for the packet to be sent. By using select() the polling and busy waiting is eliminated that would otherwise be essential for network I/O.

5.4.3(C) Clocks and Timers

The hardware clocks and timers available on most of the computer performs three basic functions:

- c) Provide the current time.
- d) Set the time.
- e) Start a cause operation X at time T.

The operating system makes heavy use of these functions. The time sensitive applications also make heavy use of these functions.

Unfortunately, the standardization of the system calls that implement these functions is not done across operating systems.

The hardware to determine elapsed time and to trigger operations is called programmable interval timer. This can be set to wait a definite amount of time and then cause an interrupt.

The setting can be for generating interrupt once or to repeat the process to cause periodic interrupts. When process is at the end of its time slice, it should be preempted.

The scheduler uses above method to generate an interrupt to preempt such process. The disk subsystem uses it to invoke the periodic flushing of dirty cache buffers to disk, and the network subsystem uses it to cancel transmissions that are proceeding too slowly because of network congestion.

The operating system may also provide "hardware timers" for user processes to use timers.

5.4.3(D) Blocking and Non-blocking VO

Q. Differentiate between blocking and nonblocking VO.

Toperaang System (MU -

IT
nace of I/O

- o. Explain issues related to performs system.

data are available. This can be full number of bytes requested. ^{er, or none at all. On the other hand} ~~and~~ ^{cal} requests a transfer that will be asynchrouous its entirety but will complete at some point in time.

- The execution of the application gets suspended after it issues a blocking system call. While executing the application was in run queue. It will move to wait queue.
- After the system call finishes, the application is moved back to the run queue, where it resumes execution.
- When it resumes execution, it will receive the values returned by the system call. The physical actions performed by VO devices are usually asynchronous. They take a varying or unpredictable amount of time.
- However, most operating systems use blocking system calls for the application interface. This is because blocking application code is easier to understand than non-blocking application code.
- Apart from blocking VO, some of the user-level processes require non-blocking I/O. One example is a user interface that receives keyboard and mouse input while processing and displaying data on the screen.
- Overlapping of execution with VO can be achieved by writing multithreaded application. While some threads carry out blocking system calls, at the same time others carry on execution. The Solaris developers used this technique to implement a user-level library for asynchronous I/O, freeing the application writer from that task. Some operating systems provide non-blocking VO system calls.
- If the call is non-blocking, then it does not stop the execution of the program for an extensive time. It returns promptly, with a return value that shows number of bytes transferred.
- The substitute to a non-blocking system call is an asynchronous system call. An asynchronous call returns immediately, without waiting for the VO to complete. The application continues to execute its code.
- The completion of the VO at some future time is communicated to the application, either through the setting of some variable in the address space of the application or through the triggering of a signal or software interrupt or a call-back routine that is executed outside the linear control flow of the application.
- As compared to asynchronous system calls a non-blocking readQ returns without delay with whatsoever

Topic : K.mel I/O Subsystem

5.4.4 Kernel VO Services

- Q. Explain various services provided by kernel are related to I/O.

Following services are provided by kernel which are related to I/O.

Services provided by kernel related to I/O.

- (A) I/O Scheduling
- (B) Buffering
- (C) Caching
- (D) Spooling and Device reservation
- (E) Error handling

Fig. C5.13 : Services provided by kernel

→ 5.4.4(A) I/O Scheduling

- It involves deciding the good order in which I/O requests should be executed. Scheduling can enhance the overall system performance. Processes also share the device access in a fair manner due to scheduling. The average waiting time for I/O to complete also can be minimized.
- There is a wait queue of requests maintained for each device. The I/O request of the application is placed on the wait queue for that device when application issues blocking I/O system call. All these requests are then arranged in a manner such that, the overall system efficiency and the average response time experienced by applications will be improved.
- Operating system tries to give equal service to all the applications. Delay-sensitive requests may be serviced on the basis of priority. For asynchronous I/O, kernel should keep track of many I/O requests simultaneously. Hence, the OS might attach the wait queue to a device's status table. This table keeps the entry for each VO device and its management is carried out by kernel.

table entry shows the device's type, address, and state of the device can be noted.

** <>**!■
if device is busy in process then
if * tyt* and o, hCr parameters will be stored in the
then device. Sched, the operations improves performance
If it is in the finance SVSt, Ott.
< in main memory or on disk via teeh?* storage
Peril's caching, and PooUng.

4.4(0) Buffering

First reason for buffering is to handle data transfer between the producer and consumer of data stream. The second reason for buffering is different data transfer sizes of the different devices. Such similarities are particularly common in computer networking. Where fragmentation and reassembly usages buffers are heavily used. Buffering is also done to support copy semantics for applications. Various buffering needs to be considered for block and character devices for a diversity of reasons. For example when process reads data from a modem, user process does a **read system call and block waiting for one character**.

Each incoming character causes an interrupt. The interrupt service procedure gives the character to the user process and unblocks it. After pulling the character somewhere, the process reads a further character and blocks another time.

- The problem with this approach is that the user process needs to be started up for every arriving character. Permitting a process to run again and again for short runs is inefficient, so this design is not a good one.
- An enhancement on above approach is to have user process an n -character buffer in user space and does a read of n characters. The interrupt service procedure places arriving characters in this buffer until it fills up.
- After this, it wakes up the user process. This proposal seems to be more efficient than the previous one. The limitation of this scheme is that if buffer is paged out when a character arrives, the buffer could be locked in memory. In this way, if many processes initiate locking pages in memory, the pool of available pages will reduce in size and performance will degrade.
- In other approach, a buffer can be created within the characters buffer, the page

Storage Management

with the user buffer is known as the buffer cache. This approach arrives while brought in from the main memory. As the buffer over this is to have a fort I has em * 7 first buffer. Within the second buffer copied to the second buffer is being copied to user space, like this. In this way, the two buffers can be used for buffering. Also important on output.

5.4.4(C) Caching

The cached compared to original, for example, disk stores instructions of currently executing process. Same instructions are then cached in physical memory, and copied again in the CPU's secondary and primary caches.

Buffer holds the only existing copy of a data item, whereas, cache holds a copy on faster storage of an data item that resides in another place. Sometimes, a part of memory area can be used for both buffering and caching. For efficient scheduling of disk VO, operating system uses buffers in primary memory to hold disk data. The VO efficiency for files shared by many applications is improved by using these buffer maintained in main memory as a cache. The files being written or reread speedily also makes use of these buffer as cache.

For the file VO request, kernel first checks buffer cache in main memory to avoid physical disk VO. In addition, disk writes are collected in the buffer cache for several seconds. This helps to collect the large transfers to **permit efficient write schedules**.

5.4.4(D) Spooling and Device Reservation

Some devices, such as printer cannot accept interleaved streams. A buffer that holds output for a device. Printer can print a single job at a time.

X «

want to print the output

simultaneously, then each one's output is spooled to a separate disk file. Once the printing is done, the spooling system queued the corresponding spool or output to the printer. queued spool files then are copied by spools system to the poorer system. In other operating system, Zme openmag system is managed by enm-kent thread.

For this purpose, opera "J d y SBNia dminitire K.R. interface that allows users to eliminate unwanted jobs prior to display the queue, to while the printer is those jobs print, to suspend pnm g serviced, and so on.

Devices like tape drives and printers, cannot helpfully multiplex the I/O requests of multiple parallel applications. Spooling is one means operating systems can manage concurrent output. Other method to handle concurrent device access is to offer explicit facilities for coordination.

■4.5.4.4(E) Error Handling

- If an operating system uses protected memory, it can protect against many types of hardware application errors. Hence, complete system failure is avoided if minor failure occurs. There can be a transient reason to fail devices and I/O transfers due to network becomes congested. The permanent reason for such failure can be, for example disk controller becomes faulty.
- In case of transient failure, operations can be repeated to recover from failure. But it is not possible to recover in case of permanent failure.
- An I/O system call returns single bit as a status indicating success or failure. UNTX uses integer variable named error to return an error code. Nearly about hundred values are used to indicate the nature of the failure.

5.4.5 I/O Protection

- User process should not issue any illegal I/O instruction to disturb the operation of system. To prevent the users from performing illegal I/O operations, all I/O instructions are defined as privileged instructions.
- User program has to execute system call in order to

oZ' ne r a n 8 SyStem Cany 0Ut 1/0 On behalf
ehZ JJeth 8 SyStem nning " m o nitOT mode
valid, then carry out trie I/O requested.

Memo* protection fivs tem should protect any memory locations from mapped and I/O po~~ss~~ memory locations. Kernel cannot just reject all user access.

5.4.6 Kernel TMta Structures

Kernel keeps state information regarding the use of I/O through different in-kernel data structures. It compiles openLS file table structures, kernel uses just like the many structures for different VO activities like network connections, etc.

Syllabus Topic: Ira „itJiin*nil 1/0 Requests to Hardware Operations

Transforming I/O Requests to

5.4.5 Hardware Operations

Q. Explain i@ - RWU es , is transfo, m ,o hard *are operations.

This plainS the t^o W o S C o D neCtS " P₁ 08ram req " e S , set of network wires or to a specific disk in order to read a file from disk, the application

JX W <he data by a file name. Instead of the disk, this file name is mapped by file system through the file-system directories to get the space allocation of the file.

In MS-DOS, file name's first part preceding the colon, for example C: indicate that it is a primary hard disk. Then, c: is mapped to exact port address through a device table. Colon separates the device name space from the file-system name space within each device.

Modern operating systems get important flexibility from the several stages of lookup tables in the path between a request and a physical device controller. The methods to pass requests between applications and drivers are common. Hence, new devices and drivers can be added without recompiling the kernel. Some operating systems can load device drivers on demand.

Following are the steps involved in typical life cycle of a blocking read request. It shows I/O operations consumes many CPU cycles.

- A process makes a blocking read () system call to a file descriptor of a file that has been opened previously.
- Kernel code related to system call verifies the parameters for correctness. If required data is in the buffer cache, then it is returned to the process which completes the I/O operation.
- If required data is not in buffer cache then a physical I/O must be carried out. The requesting process is then

removed from the run queue to the wait queue of the device, and the I/O request is scheduled. Finally, the I/O subsystem sends the request to the device driver. The request is either sent via a subroutine call or an in-kernel message depending on operating system.

4. The kernel buffer space is allocated by device driver to take delivery of the data and schedules the I/O. Finally, the driver issues commands to the device controller by writing into the device-control registers. The device controller controls the device hardware to carry out the data transfer.

5. The status and data may be polled by driver, or it may have initiated a DMA transfer into kernel memory. DMA controller manages the transfer, which generates interrupt after completion of transfer.

6. The wa ct interrupt handler gets the interrupt through the intenr.pt vector table. It signals the dev.tee driver, and returns essential data from the TW device driver gets the signal. The request has finished, finds the status of "0" and signals the kernel I/O subsystem that the request has been accomplished.

7. The address space of the requesting process, data or code from kernel and process gets transferred from wait queue to ready queue by kernel.

8. As process is in ready queue, it is now ready. When process gets CPU, it resumes execution at completion of the system call.

Syllabus Topic : STREAMS

54.6 STREAMS

Q. What is STREAMS? Explain.

STREAMS are a mechanism offered by UNIX V system. This mechanism allows an application to assemble pipelines of driver code dynamically. A stream is a full-duplex link between a device driver and a user-level process. It consists of a stream head, a driver end, and zero or more stream modules. Stream head interfaces with user processes. The device is controlled by driver end. Zero or more stream modules exist between stream head and driver end as shown in Fig. 5.4.2.

The read queue and write queue is associated with each module, stream head and driver end. The functionality of STREAMS processing is as follows. The ioctl() system call is used to push data onto a streams.

Storage Management
The exchanges of messages take place between queues in adjacent modules. Therefore, queue present in one module may overflow an adjacent queue. Hence queue should support for flow control.

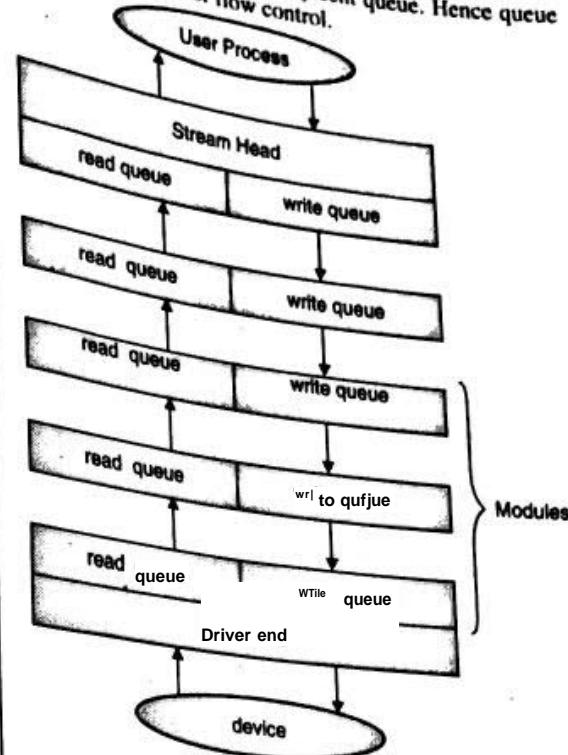


Fig. 5.4.2 : The STREAM structure

If sufficient support is available, messages are exchanged among queues in adjacent modules.

`write()` or `putmsg()` system call is executed. The `write()` system call is executed to write raw data to the stream, while `putmsg()` permits the user process to state a message. For any of this system call the stream head copies the data into a message and sends it to the queue for the next module in line.

This copying of messages carries on till the message is copied to the driver end and so the device. In the same way, the `read()` or `getmsg()` system call is executed by process to read data from the stream head.

STREAMS I/O is asynchronous (or nonblocking) excluding when the user process communicates with the stream head. The driver end also has a read and write queue. It is necessary to respond to interrupts. for example, one triggered when a frame is ready to be read from a network.

**Syllabus Topic : Performance****5.4.7 Performance**

- The performance of the system is affected by the I/O load hence it is important factor to consider in performance. CPU has to execute device driver code and also it has to carry out scheduling of the processes in fair and efficient manner. Also due to I/O, context switches occurs and hence, it burden the CPU and its hardware caches.
- If there are some inefficiencies in the interrupt handling mechanisms in the kernel then, I/O can expose them. Also due to I/O, memory bus gets heavily loaded during data copy between controllers and physical memory. It also happens when copy takes place between kernel buffers and application data space.
- In interrupt handling, it is required carrying out several operations. Hence interrupt handling is expensive task. During completion of I/O, process gets blocked. As a result of this blocking, context switching is required which lead to certain overhead on system performance. Several context switches and state switches takes place when there is communication between the machines in network. For example, during remote login many context and state switches takes place.

Following are the measures to improve I/O efficiency.

1. It is required to lessen the number of context switches.
2. The frequency of copying the data in memory while passing between device and application should be reduced.
3. Use large transfers, smart controllers, and polling to lessen the frequency of interrupts.
4. Use DMA-knowledgeable controllers or channels to relieve the CPU from simple data transfer in order to improve the concurrency.
5. Processing primitives should be in hardware to permit their operation in device controllers in parallel with CPU and bus operation.
6. Try to keep balance in performance related to memory, CPU, bus and I/O. If one area gets overloaded then it will cause the area to be idle.

5.5 Exam Pack (University and Review Questions)**Syllabus Topic : File System - File Control Pt**

- Q.** Explain various file attributes in brief. (Refer section 5.1.2)
- Q.** Explain various file operations in brief. (Refer section 5.1.3)
- Q.** Write short note on File Types. (Refer section 5.1.5)
- Q.** Explain the different techniques to access files. (Refer section 5.1.6)

Syllabus Topic : Access Methods

- Q.** Explain different file access methods. (Refer section 5.1.7) (10 Marks) (June 2015)
- Q.** Write short note on File Access (sequential access). (Refer section 5.1.7(1))
- Q.** Write short note on File Access (Random access). (Refer section 5.1.7(2))

Syllabus Topic : Directory Structures

- Q.** Explain different methods for defining the logical structure of a directory. (Refer section 5.1.8)
- Q.** Explain different directory operations. (Refer section 5.1.10)

Syllabus Topic : File Sharing

- Q.** Write note on "access rights". (Refer section 5.1.13)
- Q.** Explain issues related to file sharing. (Refer section 5.1.13(A))
- Q.** Explain different access rights to the file. (Refer section 5.1.13(A))
- Q.** Explain in detail operation of remote file system. (Refer section 5.1.13(B))

Syllabus Topic : Protection

- Q*** Explain file system protection in detail. (Refer section 5.1.14)

Syllabus Topic : File System Implementation

- Q.** Explain file system structure. (Refer section 5.2.1)

Syllabus Topic : Implementing File System

- Q.** Explain the implementation of file system in detail. (Refer section 5.2.2)

notes on virtual file system.
(Ref section 5.2.2(B))

Syllabus Topic : Directory Implementation
Explain approaches for implementation of directory.
(Refer section 5.2.3)

* I* <e: *kQeatlo, >W₀,

What are the different allocation
schemes to File system >s. mode with
marks)

contiguous file allocation method
'antages and disadvantages. (Refer section 5.2.4(A))

advantages and disadvantages of its
(Refer section 5.2.4(B)) (B))

Explaining indexed file allocation
disadvantages. (Refer section 5.2.4(D))

Explain the structure of i-node in detail.
(Refer section 5.2.4(D))

How i-node is used to allocate the file ?
(Refer section 5.2.4(E))

Syllabus Topic : Free Space Management

List methods of free space management.
(Refer section 5.2.5)

A Explain bit map method in detail.
(Refer section 5.2.5(A))

Explain linked list of disk blocks method.
(Refer section 5.2.5(B))

Syllabus Topic : Efficiency and Performance

Explain various techniques to improve efficiency and performance of secondary storage.
(Refer section 5.2.6)

Syllabus Topic : Recovery

Explain file system recovery in detail.
(Refer section 5.2.7)

Whjs Topic : NFS

Explain operation of network file system in detail.
(Refer section 5.2.8(A))

Explain the layered structure of NFS.

*wsecfibn 5.2.8(C))

Syllabus Topic : Storage Management

Storage Management
Physical Storage

Explain physical structure of magnetic disk.
(Refer section 5.3.1)

Explain physical structure of magnetic tapes.
(Refer section 5.3.2)

Syllabus Topic : Storage Management

Explaining inhost storage
(Refer section 5.3.3(A))

Explain network.

Explaining disk scheduling

• JZ " " scheduling
(Refer section 5.3.4(A) to 5.3.4(E)) (10 Marks)

Example 5.3.3 (10 Marks) (Dk. 2016)

Example 5.3.10 (10 Marks) (Dec. 2014)

Example 5.3.11 (10 Marks) (May 2016)

Example 5.3.12 (10 Marks) (June 2015)

Example 5.3.13 (10 Marks) (Dec. 2015)

Syllabus Topic : RAID Structure

Q. Explain various RAID levels.
(Refer section 5.3.7(A))

Syllabus Topic : Tertiary-Storage Structure

Q. Explain various Tertiary-Storage Devices.
(Refer section 5.3.9(A))

Syllabus Topic : Swap-Space Management

Q. Explain management of the swap space in detail.
(Refer section 5.3.10)

Syllabus Topic : I/O systems

Q. Explain different types of I/O devices.
(Refer section 5.4.1(A))

Q. Explain different parameters to distinguish I/O devices.
(Refer section 5.4.1(B))

Syllabus Topic : Overview of Hardware

Q. Write note on "Polling". (Ppts' on 5.4.2(B))

- Q. What is interrupt? Explain the tasks carried out by interrupt handler. (Refer section 5.4.2(C))
- Q. Explain steps in DMA transfer. (Refer section 5.4.2(E))
- ☛ Syllabus Topic : Application I/O Interface
- Q. Differentiate between blocking and nonblocking I/O. (Refer section 5.4.3(D))
- Q. Explain issues related to performance of I/O system. (Refer section 5.4.3(D))

- ☛ Syllabus Topic : Kernel I/O Subsystem
- Q. Explain various services provided by kernel related to I/O. (Refer section 5.4.4)
- ☛ Syllabus Topic : Transforming I/O Requests to Hardware Operations
- Q. Explain I/O Request is transformed to hardware operations. (Refer section 5.4.5)
- ☛ Syllabus Topic : STREAMS
- Q. What is STREAMS? Explain. (Refer section 5.4.6)

CHAPTER

6

Distributed Systems

Module VI

Syllabus Topics

Distributed Operating System : Network based OS ,
Structure and Protocols; Distributed Service, File p-
trof and Deadlock Handling, Naming and distributed
Svnrk and To Po lo gy. Communication
y, Remote file access, Stateful
ynterntzahon : Mutual Exclusion, Concurrency

Syllabus Topic : Distributed Operating System

Distributed Systems

Q. What is a distributed system.

The development of powerful microprocessors and invention of the high speed networks are the two major developments in computer technology.

Many machines in the same organization can be connected together through local area network and information can be transferred between machines in a very small amount of time.

As a result of these developments, it became easy and feasible to organize computing system comprising large number of machines connected by high speed networks.

Over the period of last thirty years, the price of microprocessors and communications technology has constantly reduced in real terms. Because of radical distributed computer systems appeared as a practical substitute to uniprocessor and centralized systems.

The networks of computers are present all over. Internet is composed of many networks which share networks separately and in a non-pertinent manner. The necessary characteristics that make topics to focus under distributed system

6.1 Definition

A computer network is defined as a set of communicating devices interconnected by

communication links. These devices include computers, printers and other devices capable of sending and/or receiving information from other devices on the network. These devices often called as node in the network. So computer network is interconnected set of autonomous computers.

- A distributed system is defined as set of autonomous computers that appears to its users as a single coherent system.
- Users of distributed system feel that, they are working with a single system.

Main characteristics of distributed system

Q. What are its characteristics? Explain.

Following are the main characteristics of distributed system,

- A distributed system comprises computers with distinct architecture and data representation. These dissimilarities and the ways all these machines communicate are hidden from users.
- The manner in which distributed system is organized internally is also hidden from the users of the distributed system.

The interaction of users and applications with distributed system is in consistent spite of where and when interaction occurs.

- A distributed system should allow for scaling it support for availability. It should be reliable to the users and applications of failures.

X handling should be hidden and application

6.1.2 Motivation

Q. What are objectives behind building the distributed system?

Following are the objectives behind building distributed system.

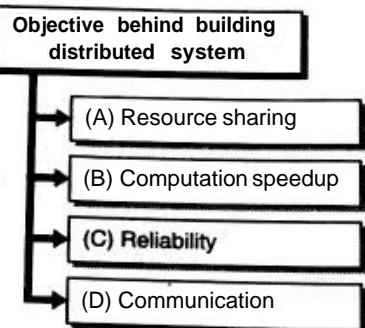


Fig. C6.1 : Objective behind building distributed system

→ 6.1.2(A) Resource Sharing

- Resource sharing offers saving in cost. One printer can be shared among many users in office instead of having one printer to each individual user. There can be more saving in cost if expensive resources are shared.
- The increase in connectivity and sharing also increases the security risk and to deal with it is equally important. Presently, systems offer fewer defenses against eavesdropping or intrusion on communication.
- A communication can be tracked to construct a favorite profile of a particular user. This clearly violates privacy, particularly if it is done without informing the user. A allied problem with increased connectivity can also cause unnecessary communication, for example electronic junk mail, called as spam. Special information filters can be used to select inward messages based on their content.

→ 6.1.2(B) Computation Speedup

- Many big computations can be divided in sub computations which can be assigned to different nodes in network so that work can be carried out in parallel.
- In this case, speed of computation increases. As another example, process may be migrated from heavily loaded machine to lightly loaded machine in network.

→ 6.1.2(C) Reliability

- System does not affect with failure of one node in network as other nodes are available in system.

- In case of failure, other nodes can take over computations to achieve reliability. This failure should be detected and recovered by system itself. Once the recovery is done, final computation must result in consistent state.

→ 6.1.2(D) Communication

Machines in network communicate with each other by exchanging messages. At low level, this communication takes place through protocol stack between machines.

At highest level, applications exchange these messages. Due to such communication, many users can be part of a single work to be carried out by sharing the information.

- By connecting users and resources, it becomes easier to work together and exchange information. The success of the Internet is due to its straightforward protocols for exchanging files, mail, documents, audio, and video.
- The worldwide spread people can work together by means of groupware. Electronic commerce permits us to purchase and sell variety of goods without going to shop or even leaving home.

6.2 Types of Distributed Operating Systems

Q. Explain different types of distributed operating systems.

- In a network, multiple machines may have different operating system installed on it. Operating systems for distributed computers are categorized as :

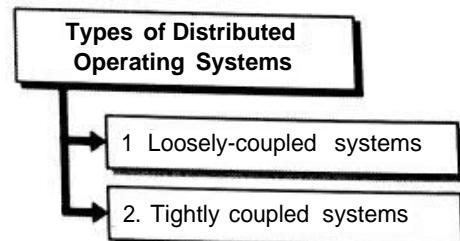


Fig. C6.2 : Types of Distributed Operating Systems

→ 1. Loosely-coupled systems

In a set of computers, each has its own OS and there is coordination between operating systems to make their own services and resources available to the others. This loosely coupled OS is called as **network operating**

Tightly coupled systems

If TCP/IP is a single global view of the network, such that it is coupled, called distributed operating system. It is useful for the management of multicomputer systems and underlying hardware. This hardware consists of many processes and details that are hidden.

Z Syllabus Topic : Network Based OS

6.2.1 Network Operating System (NOS)

Users have NOS installed on machine A. They can access resources on other machines in the network by logging in or transferring data from a remote machine to their own machine.

6.2.1(A) Remote Login

- Network operating system allows users to log in remotely. The Internet offers the telnet facility for remote login. NOS permit users to use services available on a particular machine in the network. Remote login service provided by NOS allows the user to log in remote machine from his/her terminal.
- Using command for remote copy, user can copy the file from one machine to other.

6.2.1(B) Remote File Transfer

- NOS offer a mechanism for remote file transfer from one computer to another. Here, each computer has its own local file system. If a user working on computer A wants to access a file on another computer B, then the file must be copied explicitly from the computer B to computer A. File Transfer Protocol (FTP) is used for such transfer in internet environment. User invokes the FTP program as:

ftp name of computer B

After entering the correct username and password, the user should connect to the right subdirectory to access the required file. Suppose user wants to access the file abc.frt from computer B. Then this file must be copied to computer A by executing following command:

get abc.txt

The location of distributed system is transparent to user. i.e. "copy from local to remote" is same as "copy from A to B".

Whichever is the space, the user connects to it and performs operations. This completes the task.

b. get: Fetch file from remote computer, computer to local.

2. Put: transfer file from local computer to remote.

3. Wdir: Usiflik is in the current directory on the computer.

4. cd: Change current directory to remote.

6.2.2 Distributed Operating Systems (DOS)

local resources are to be the same way. Operating systems use data and process migration from one site to another.

- (A) Data Migration
- (B) Computation Migration
- <S process Migration

6.2.2(A) Data Migration

- User working on machine A can access data (such as a file) that reside at machine B. One method of data migration is to transfer the complete file to machine A. After this access to the file is local. When the user needs a copy of the file if it has been modified is sent back to machine B.
- Although a small change has been made to a large file, all the data needs to be transferred. This method was used in the Andrew file system. This method of file transfer is inefficient. In other method of data transfer only those portions of the file that are actually needed for the immediate task gets transferred.
- If another part of the file is needed later on, another transfer will take place. When the user no longer wants to access the file, any part of it that has been modified is sent back to machine B. The Sun Xylems Network File System (NFS) protocol and Microsoft SMB protocol are used on top of either TCP/IP or the Microsoft NetBEUI protocol.

SeS (running) 3150 aows file sharing over a network.

6.2.2(B) Computation Migration

- The computation rather than the data can be transferred across the system. This approach is called computation migration.
- If applications need to access various large files that reside on different machines, to obtain a summary of those files. In this case, it would be more beneficial and efficient to access the files at the machine where they reside and return needed results to the machine that initiated the computation.
- If the time required to transfer the data is longer than the time to execute the remote command, the remote command should be executed.

6.2.2(C) Process Migration

- Q. What are the reasons behind process migration?

Process migration involves transfer of process to other machine for execution purpose. The whole process, or portion of it, may be executed at different machines. The reason for migration is :

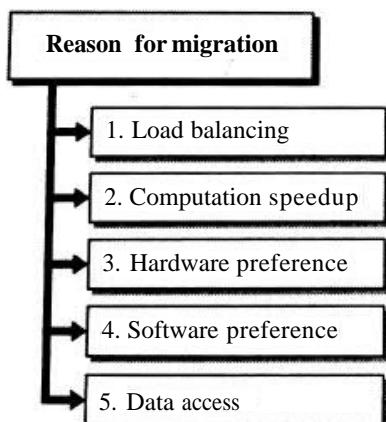


Fig. C6.3 : Reason for migration

→ 1. Load balancing

The processes (or sub processes) may be distributed across the network on different machines to balance the workload.

→ 2. Computation speedup

The subprocesses of the single process can run concurrently on different machines then the total process turnaround time can be reduced.

→ 3. Hardware preference

The particular process may need some specialized processor for execution.

→ 4. Software preference

The needed software by process may available at only a particular machine, and either the software cannot be moved, or it is less expensive to move the process.

→ 5. Data access

If the huge amount of data required in the computation it may be more efficient to have a process run remotely than to transfer all the data.

Syllabus Topic : Network Structure

6.3 Network Structure

- Distributed systems are built on top of computer networks, which are of two types: LANs (Local Area Networks) and WANs (Wide Area Networks).
- LAN covers one room, building or campus, WAN covers a city, country, or whole world. Ethernet is an example of LAN.
- Internet is a collection of thousands of separate networks and can be considered as one WAN.

6.3.1 Local-Area Networks (LANs)

- Q. Explain LAN type of network.

- Instead of having large mainframe computer, it seemed to be more economical to have many several computers. Hence, LANs emerged in early 1970 as a substitute to large mainframe systems which many organizations were using for their computation.
- As LANs cover small area such as LAN covers one room, building or campus, the communication links between nodes are having more speed and less error rate.
- The reliability and high speed can be achieved by using high quality cables to connect the different computers in network.
- Twisted pair cables or fiber optics cables are most commonly used in local area network.
- IEEE Standard 802.3 describes classic Ethernet in which a coaxial cable is used to connect many computers. The cable is called the Ethernet. In the very first version of Ethernet, a **vampire tap** was used to attach a computer to the cable. Connecting the machine to cable was carried out by drilling a hole halfway through the cable and screwing in a wire leading to the computer.

Multiaccess bus, ring, and star networks are the most configurations for local area network.

Very few networks and Blu, AX work have communication speed of 1 Gbps. Radio per second.

Gigabit Ethernet have speed of 1 gigabit per second. 10BaseT Ethernet have speed of Ten megabits per second. 100BaseT Ethernet requires a higher-quality but runs at 100 megabits per second and is becoming common.

contains several small and large nodes which peripheral devices which are shamable, and fP to connect to other networks.

Wide-Area Networks (WANs)

Explain WAN type of network.

WANs cover a city, country, or the whole world. A met collection of thousands of separate networks can be considered as one WAN. WAN was 30 in late 1960 as an academic research project to offer efficient communication among machines. The main objective behind this project was to share hardware and software in convenient and economical manner by a wide community of users. Arpanet is the example of first WAN developed.

WAN machines are physically distributed over a geographical area. Therefore communication are relatively slower compared to LAN and links are also not reliable. Internet WAN offer capability for communication between different machines which are geographically distributed.

The Internet comprises hosts and routers which are also

computers. Hosts are PCs, notebooks, handhelds, mainframes, and other computers owned by individuals or companies that want to connect to Internet.

Routers are also computers with many incoming and outgoing lines. It accept incoming packets on their incoming lines, process it and end them on their way along one of many outgoing lines. Routers are connected together in large networks.

Large national or worldwide route and ISPs (Internet Service Providers) for their customers. Internet is organized with backbones at the top. Backbone operator. Backbone central

Distributed Systems
faceted by high bandwidth fiber optics. Telephone companies have connections to backbone.

Host are not directly connected to backbone. Regional backbone by routers are connected to routers in Corporate backbone.

Network routers are connected to regional backbone. Users at ISPs are connected to every host. USed IS Customer in this way, often on internet has at least one path, and to every other host.

Syllabus Topic : Network Topology

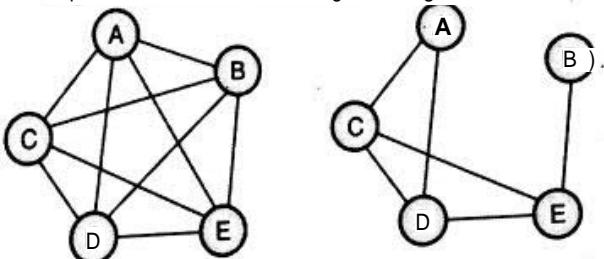
Network Topology

Q Explain different topology for network with its advantages and disadvantages.

Machines in distributed system can be connected in different ways called as topology for that network. Each topology has its advantages and disadvantages.

Installation cost, communication cost and availability criteria are used to compare different configurations.

Various topologies are shown in Figs. 6.4.1, 6.4.2, 6.4.3 and 6.4.4. In fully connected networks, each node is connected to every other node in network. Hence, number of links increases as square of number of sites. Also installation cost is high. Hence, this topology is impractical to consider building the large network.



(a) Fully Connected Network (b) Partially Connected Network

Fig. 6.4. i

- In partially connected networks, there exist direct link between nodes, but not between all nodes as in fully connected network. Therefore installation cost is less than fully connected network. As there is no direct link between nodes, the communication cost is higher. Because of failure of some links, some of the nodes may become unreachable.
- The structured networks, ring and star networks are examples of partially connected networks.

networks. Tree-structured network relatively have low installation and communication cost. However, due to the failure of a single link, network may become partitioned.

- In ring network, if two links fails then the network can become partitioned. So availability of ring network is higher than tree-structured network. Browsing network incur high communication cost as messages needed to cross large number of links to reach to destination.
- In star network, if single link fails, then only single node will be unreachable. In star network, communication cost is low but if central node fails, then all nodes in the system become disconnected.

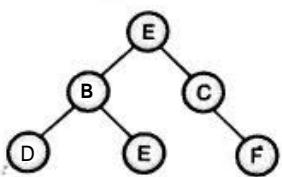


Fig. 6.4J: Tree-structured Network

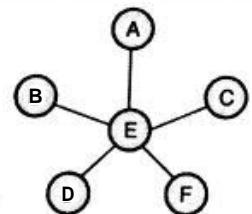


Fig. 6.43: Star Network

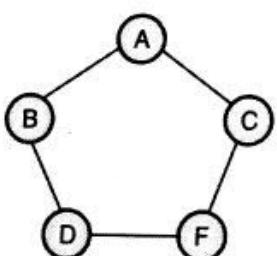


Fig. 6.44 : Ring Network

Topic: Communication Structure

6.5 Communication Structure

Following five different issues are important to address while designing the communication network.

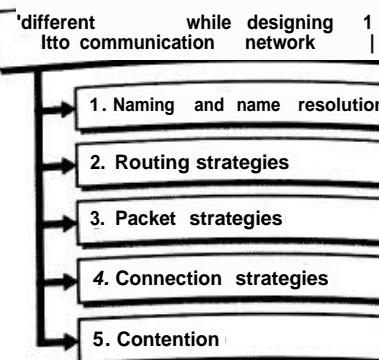


Fig. C6.4 : Issues while designing communication network

-4 6.5.1 Naming and Name Resolution

Q. what is naming? How names are resolved? Explain.

- Processes running on different machines must be able to specify each other. On each machine, process has process identifier. The messages to be sent to the processes may be addressed to the process identifier. As machines do not share memory, initially machines have no knowledge about processes on remote machines.
- To solve this problem, a pair <host name, identifier*> is used to identify the process on remote machine. Host name is unique in network and is alphanumeric so that it is easy to recognize it. Within any host, identifier is either process identifier or other unique number used in that host to identify the process.
- Names are used to identify the hosts as they are easy to remember. But machines (hosts) use an identifier for simplicity and speedup. Hence, it is necessary to have some mechanism to resolve host name into host*id. Two solutions exist. First each host may have data file which contains host name and addresses of all hosts in network. This was the approach used initially in internet. But this approach requires updating the data files on all machines whenever hosts are added or removed from network. Due to growth of internet, this approach is not used for name resolution.
- The second approach is domain-name system (DNS) which specifies the naming structure of the hosts, and also name-to-address resolution. Names start with most specific part, where each part is separated by dot operator and ends with general part. For example, *aihu.ityale.edu* indicates host name *aihu* in information

technology department at yale - 6-7
domain.edu. Uln leisity in the
address each component h_f
, ss io system) which takes a name server
C the "" T" rcSPOn * <>le I?* the
Finally, name server for the host is .^{that} ante.
then host-id. Consider Machine J, 6, 1 which
X t is , issued by process X on which
request athu.it.yale.edu. Following steps t the host
resolve the name athu.it.yale.edu. lnv oived to

Kernel on host X issue request to name server of edu domain. It ask to name server about name. The edu name server's address UV of interface initial request was issued. know >

jjK?w -----ne server r----- "duress Of yale.edu"

Kernel then request to name server v,i, asking at it-pccrtf it x- j j.u.

To the return address by it.yale.edu kernel issues request for address of athu.it.yale.edu. Finally IP address (host-id) of athu is returned.

465-² Routing Strategies

Q. Explain different routing strategies.

Machines in internet communicates by exchanging packets. Each packet contains address of source and destination. Router contains routing tables. It extracts address in incoming packet and look up the table to find which outgoing line to send the packet on and thus to which router.

- This process is repeated until the packet is delivered to the destination host. The routing tables are very dynamic and are updated continuously as routers and links go down and come back up and as conditions change.

Following are the three most common routing methods used for routing purpose.

Three most common routing methods used for routing purpose

- 1 Fixed routing
 - 2. Virtual routing
 - 3. Dynamic routing

Flag.C6 : Rou « W m£lh » dS

Distributed Systems

1. Fixed routing

In this type of routing, path from sender machine to destination machine is specified in advance. This path changes if and only if a hardware failure disables it. Otherwise, it does not change.

2. Virtual routing

In this type of vestigation ntJTM¹¹⁸, P>fill f_u>ra "■>< machine to S₈s₈®lon fixed for the duration of one

Different paths --- Used the ~~message~~ f y different sessions to send U OmSOURct destination.

** 3, Dynamic routing

In dynamic routing, the path used to send a message **source to destination machine is selected only when a message is sent.** As decision is taken dynamically, different paths are used by different messages.

- in general, a source machine sends a message to destination machine on whatever link is the least used at that particular time.

→ 65.3 Packet Strategies

- There are two types of services : Connection-oriented service and connectionless service. With connection-oriented network service, user initially establishes a connection, uses the connection, and then releases it. Over this connection, the sender sends bits at one end, and the receiver takes them out in the same order at the other end.

- On the contrary, with connectionless service packets, routes to reach the same destination are out of order at

r . " " -

oriented service.

- Services are characterized by a quality of service.
- Reliable service is implemented by giving a guarantee to sender that packet is correctly delivered to receiver.
- Receiver SCM's Acknowledgement packet as a confirmation of delivery of packet to the sender.
- "cohfiraaho" of acknowledgements introduce delay and are essential to detect packet loss.

- slow the communication. File transfer is example of reliable connection-oriented service.
- Reliable connection-oriented service supports either message sequences or byte streams.
- In message sequences message boundaries are preserved, in byte streams, the connection is simply a stream of bytes, with no message boundaries.
- For application such as digitized voice traffic, the delays introduced by acknowledgements are unacceptable. Here noise on telephone line or a garbled word from time to time user can prefer with compare to delay due to waiting for acknowledgements.
- Unreliable connectionless service is called as datagram service and does not offer acknowledgement as a confirmation of delivery of packet to the sender.
- The acknowledged datagram service can be used to send one short message instead of establishing the connection. The request-reply service can be used for client server communication. Following are the six types of network services.

*+ 6.5.4 Connection Strategies

Q. What are the different connection strategies used in network communication?

- Once connection is established between source and destinations, processes can set up communications sessions to exchange information. There are number of ways by which two communicating processes may connect. Following are the three most common schemes used.

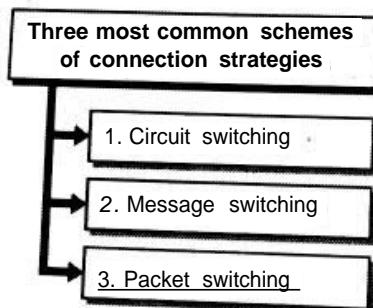


Fig- C6.6 : Schemes of connection strategies

→ 1. Circuit switching

- system, connection established in telephone
- Two communicating parties continue to talk over a communication line till connection is terminated. During communication period, no one can use this line

- A permanent physical link is established between the two processes that want to communicate with each other.
- Any other process cannot use this established link till communication session is completed. Hence in circuit switching, link is allocated to for the duration of communication session.

→ 2. Message switching

- In message passing, a temporary link is established for the duration of one message transfer between two communicating processes.
- As per need, physical links are allocated dynamically among correspondents. These links are allocated for only short duration.
- Each message contains data with other system information like source, the destination, and Error Correction Codes (ECC). This information is used by the network for correct delivery of the message to the destination.

→ 3. Packet switching

- Each message from application is first broken up in small chunks called a packet.
- Each packet contains data with source and destination address as each packet may be sent to its destination separately. Each packet reaches to the destination machine through separate path.
- As packets takes different route to reach the destination, all must be reassembled into messages as they arrive.

→ 6.5.5 Contention

If several computers want to transmit over a same communication link then collision occurs. This situation occurs in ring and multi-access bus network. To avoid the collision as it degrades performance, several techniques have been developed.

CSMA/CD

- In this technique, every computer sense the channel before transmitting to check that other machine message is currently being transmitted or not. This is called carrier sense multiple access technique. If collision is detected, then transmission gets terminated.
- On Ethernet, computer has to first see whether other computer is transmitting packet or not. Then it can send the packet. If two or more computer sends packet at the

same time collision occur and both detect it. So until Lijst computer fini shes sending, other computer can not sebdP et.

Vision occurs, both the issions. Here both computers wait for amount of time between 0 and T Microsecond. In 6-9 vision occurs, all colliding vision occurs, all colliding

ji into the interval 0 to $2T$ M again. On each collision T^6 maximum Visions. This algorithm is uponed⁸¹ backoff.

Si set on Ethernet for also a maximum number of cable connected to it, a multiple of computers to be entire campus.

A bridge is used to connect these which permits traffic from another when the source and destin

Ethernets use switches to avoid collisions. Each switch contains ports, to which computer another switch can be attached. A bridge is connected to a switch and sent out on the port where destination machine stays. As each machine is connected to separate port and collision is eliminated but requires bigger switches to connect many computers.

Token Passing

Token is a unique message that is circulated in network. Any machine want to transmit holds the token. Once transmission completes, it releases token so that other computer can use it. Token may loss so system should be able to generate new token.

Syllabus Topic : Communication Protocols

6.6 -Communication Protocols

Explain different communication protocols.

Protocol is set of rules by which computer communicates with each other. Many protocols exist such as router-router protocols, host-host protocols, and others. Protocol stack layers protocols on top of one another. This protocol stack is

used by all modern systems. Different layers in protocol stack deal with different issues.

International "Hid" organization (ISO)

Following are the functions of each layer during communication between two machines.

Physical layer: This layer deal with mechanical, electrical, functional and process characteristics of bits per second.

1 5 H physical UTM, "1 many volts to represent 0 and 1. These issues are simple, whether these issues are defined by physical layer. It is implemented in hardware of network.

Data link layer: It is a device used to connect two network nodes. It is responsible for channel allocation and transmission.

Network layer: All, k

The layer offer, for routing the pack

IP layer provides logical address and maintains

Transport layer: It is responsible for transfer of messages between clients, including breaking the messages into segments, sequencing of packets, control, and generating physical addresses.

Session layer: The session layer is responsible for establishing the sessions between users. It also implements process-to-process communication protocols. It provides services like dialog control, synchronization, token management.

Presentation layer: This layer deal with syntax and semantics of information exchanged between machines.

As machines may have different data representations, presentation layer resolve the differences in formats among the various machines in the network, including character conversions and half duplex-full duplex modes (character echoing).

Application layer: This layer interacts directly with users.

It contains different protocols needed by users such as TFTP, Telnet, DNS, and SMTP. This layer deals with file transfer, remote login protocols and electronic mail.

- Most of the distributed systems use the Internet as a base. Hence, these system uses two important Internet protocols: IP (Internet protocol) and TCP (Transmission Control Protocol), IP (Internet Protocol) is a datagram protocol. In IP, a sender sends datagram, of up to 64 KB over the network and no guarantees are given for its delivery. The datagram may be fragmented into smaller packets and travel independently, possibly along different routes. At destination host when all packets reaches, they are assembled in correct order as per sequence number and delivered to the application.
- IP protocol has two versions, v4 and v6. Version 4(V4) are currently in use and v6 is up and coming. IP v4 packet starts with a 40-byte header that contains each 32-bit source and destination address with other fields. These are called IP addresses and routing is carried out using these addresses.
- IP does not offer reliable communication in the Internet. To offer reliable communication, TCP (Transmission Control Protocol), is present on top of IP. TCP makes use of IP to offer connection-oriented streams. The remote process always listen the incoming connection on port number. Remote process is specified by IP address of machine and port number. Sender first establishes the connection and sends bytes over that connection which is guaranteed to come out the other end undamaged and in the correct order. The TCP gives this guarantee by using sequence numbers, checksums, and retransmissions of incorrectly received packets.

Syllabus Topic : Distributed File Systems

6.7 Distributed File Systems (DFS)

Q. Explain working of distributed file system (DFS).

- In distributed system, files are available on several computers. Computers in distributed system can share these physically dispersed files by using distributed file system. Service offers particular function to client and it is a software entity running on some machines. Server runs service software on single machine. Client process invokes the service through some set of defined operations called as client interface.
- File system offers file services to clients. The set of primitive file operations are create a file, delete a file, read from a file, and write to a file. Client interface is formed with set of these operations. File server controls

storage dCViCCS & UCh " on Which met are stored. These files are accessed from these devices as per request of client.

- There can be different implementation for DFS. Server may run on dedicated machines. In other implementation both client and server may run on same machine DFS can be part of DOS or it can be a software layer managing communication between file system and NOS. DFS should come into view to its client as conventional centralized file system. The servers and storage devices which are on different machines in network should be invisible to clients. DFS should fulfill the request of client by arranging the required files or data.
- As data transfer is involved in operation of DFS, its performance is measured with amount of time required to service the client request. Storage space managed by a DFS includes different and remotely located small storage spaces.

Syllabus Topic : Naming and Transparency

6.8 Naming and Transparency

- Mapping between logical object and physical object is naming. User always knows logical name of the file but system has to manipulate the data blocks which is actually stored on disk sectors or tracks. Actually text name of file gets mapped to numerical identifier which in turn again mapped to disk block. Due to such mapping user remains unknown about location of file on secondary storage.
- In case of DFS, location of file in network remains unknown to the user. This is the new dimension added to the abstraction with compare to conventional file system. In conventional file system, file remains on disk of same machine whereas in DFS location of file can be in disk of any machine in the network. File replication is supported in DFS so that several replica of file exist. By given file name then set of locations of the file is returned by mapping.

6.8.1 Naming Structures

Name mapping should support location transparency and location independence.

- 1. Location transparency :** Physical storage location of file cannot be known by using file name.

J Ctn independence , AhhZr**

location of the Chan8es , here physical storage need to change (k. Sk" 16 , clHie nt DFS system , supports S.a. ic PPir " ted by these systems. Henri noil on of dependence is immaterial for these sterms. tion independence is suptxjrted Itu Ulen files o < lo > cal data TMn " ,ers that ar p can be , specifv borage location. If on . Cached transparency is supported, the file kx , f of physical disk blocks although J 1**1 to remain hidden from users.

With location transparency, users can shar p the files simply by naming the files in to the mole manner t B® *e files are local. In this "X " , Sdarin g the storage space is cumbersome.

location independence allows sharing of

end data objects as well. Hence, it is u con ge balancing the disk utilization across the system. n th case, naming hierarchy is separated from the hierarchy of storage devices hierarchy and front the inter computer structure. Once separation between file name and storage space completes, client can now access file from remote server also.

If clients are diskless then it relies on servers to provide all files and OS kernel. This is needed as diskless client cannot use DFS code to get back the kernel. There is a special boot protocol in ROM which is invoked by client. This facilitates networking and get backs either boot code or kernel from fixed location. As soon as kernel gets copied over network and loaded, DFS makes available all other OS files.

* In current systems, client uses both local disks and remote file servers. Operating systems and networking software are stored locally. File systems holding user data and sometimes applications are stored on remote file systems.

W Naming Schemes

1. Ptain various naming schemes.

There are three approaches available for in DFS.

In first approach, combination of f ic s stein -wide name is used. This ensures a unique name. In this, host: local -name combination!

Th-

Distribhrtm Systems

d. v. Lld' 1 L PprDKh is used in network file system "•Wortcin. " yaterns. QNC+ a vendors. NFS "PPorted by many UNIX package. NFS xr po.K " 1 of i networking cok tor " s to local a- " ra " hamsm to attach remote intrudUc ed directOT V u ? s Which appear to u based on a table of moun t feature is m ounts are done on demand, names. Restructure

h thind apDrn ach, all the files in the system in a single global name structure. This file-system is covered Saturn s n's , name structure. This file-system c entional file lsom <>rpthic to the structure of a tUKes <> diffi_e i7t tem PraCt .Cally *some special files 8 bal name structure ** om <ch the 8 of single

WoLh ? approaches , 1,1051 complex is the practically most by WS structure is complex and maintain. remote detectQry to local one. hierarchy becomes more unstructured. Due to non availability of server some directories on different computers also become unavailable,

6.3 Implementation Techniques

As transparency is important, implementation should include the mapping of a file name to the associated location. It is necessary to combine files in component units instead of keeping a single file, so that resultant mapping would be manageable. It is also helpful from administration point of view.

Replication, local caching, or both can be used for improving the availability of mapping information. If location independence is supported then consistent update of this mapping is not possible. This is due to change of locations. Solution to this problem is to use low-level location-independent file identifiers.

File names which are in text form are mapped to lower-level file identifiers that specify to which component unit the file belongs. As the identifiers are location independent, its replication and caching can be carried out freely without being invalidated by migration of component units. This will result in need of a second level of mapping which maps COMPO " n, " Nations and needs a simple yet consistent update mechanism.

- Structured names are commonly used for low-level identifiers. These names are bit strings with two parts. The first part is component unit to which the file belongs. The second part indicates a file within the unit. Implementation with more than two parts is also possible. This explained technique is used in Andrew File System (AFS).

Syllabus Topic : Remote File Access

6.9 Remote File Access

- Suppose, client forward the request to server to access remote file. Naming scheme locates the server and actual data transfer between client and server is achieved through remote-service mechanism. In this mechanism, request for accesses is forwarded to server, which then performs accesses and returns (the result to user or client. This is similar to disk accesses in conventional file system.
- Caching can be used to improve performance of remote-service mechanism. In conventional file system, caching is used to reduce disk I/O. The main goal behind caching in remote-service mechanism is to reduce network traffic and disk I/O. Following are basic caching schemes in DFS.

6.9.1 Basic Caching Scheme

If data needed for operation is not cached by client then it send request to server for the same. Now, client performs the accesses on cached data which is sent by server. All the future repeated accesses to this recently cached data can be carried out locally. This will reduce additional network traffic.

Least Recently Used (LRU) algorithm can be used for replacing the cached data. Master copy of file is available on server and its part is scattered on many client machines. If copy of the file at client side modifies then its master copy at server should be updated accordingly called as cache-consistency problem. DFS caching is network virtual memory which works similar to demand-paged virtual memory having remote server as backing store.

In DFS, data cached at client side can be disk blocks or it can be entire file. Actually more data are cached by

than actual y "d so that most of the operations are carried out locally.

- If granularity of caching is large size chunks then hit ratio increases. But, if miss occurs then more data is needed to retrieve from server which increases network traffic. It in turn will increase the possibility of consistency problems as well. The network transfer unit and the RPC protocol service unit should be taken into account while selecting unit of caching. For large caches large block sizes are beneficial.

6.9.2 Cache Location

- Cache location can either be main memory or disk. If cache is kept in main memory then modifications done on cached data will be lost due to crash. If caches are kept in disk then they are reliable. No need to fetch the data during recovery as data resides on disk. Following are the advantages of main memory caches.
 - o It permits for diskless workstations
 - o Data access takes less time from main memory compared to access from disk.
 - o Performance speed up is achieved with larger and inexpensive memory which technology demands today.
 - o To speedup I/O server caches are kept in main memory. If both server caches and user caches are kept in main memory then single caching mechanism can be used for both.

Many remote-access implementations take the hybrid approach considering both, caching and remote service. In NFS, for example, the implementation is based on remote service but is improved with client- and server side memory caching for performance.

6.9.3 Cache-Update Policy

- System's performance and reliability depends on the policy used to write back updated data to the master copy of file which is available on server. Following policies are used.

Write-through policy

- This policy writes data through to server copy (disk) the moment they are placed in any cache. This policy is more reliable as little information is lost when client system crashes. The write performance is poor as each write access has to wait until the information is sent to the server.

[S Policy there is delay > - writing the master copy on server. M caching to the master copy on server. M modifications in cache and then later time.

I L. advan, a 8c < if this J dooe on n. ten f access complete in less time as writes are made to cache. Second, the data which may be over written prior

I fatten. The limitation thi; "Mate needs to crashes when data are lost Client and hence less reliable.

There are variations of delayed choice is w flush a block < soon as " Mle > One

I selected from the chenfs cache. This altern " " be good performance, but some blocks c. lc <

I dial's cache a long time before they exist in the server. A negotiation between this eh*** 11 back K-through policy is to scan the cache and be Rivals and to flush blocks that have ij " reular siare the most recent scan. So far one " mofled 0 delayed write is to write data back to the " when the file is closed. This write-on-close policy is

used in Andrew File System (AFS),

L Consistency

I Client machine always use cached data for accesses which is consistent with master copy at server. If client examines that its cached copy is out of date, then it should get the up-to-date copy of data. Following two approaches are used.

Two approaches of consistency

- 1. Client-initiated approach
- 2. Server-Initiated approach

Fig. C6.7 : Approaches of consistency

Client-Initiated approach

I client starts a validity check in which it contacts server. In this check, client checks consistency of data with the data in master copy. resulting consistency semantics is decided by a cycle of validity checking. Validity check is carried out prior to every access or on

Distributed Systems
Every access accompanied by a validity check is delayed, compared with an access served by the cache.

h? others and, checks can be started at intervals. The load on server depends on frequency of the validity check.

2. Server-initiated approach

The server maintains the records of the files or portion of files that each client access. Server must take action when it detects a potential

If two different clients in conflicting when there is no conflict, then the server updates the cache and sends a notification to the clients.

OP mode < the intend *** when a file is specified for every open mode must be

Whenever server receives that file has been opened simultaneously in read writing modes, it can take action by disabling cache. When a file is opened, it is disabled. When the operation becomes active, it is enabled.

6-9 5 A Comparison of Caching and Remote Service

Q. What is the difference between remote service and caching?		
Sr. No.	Caching	Remote Service
1.	It allows to serve remote accesses locally so that they can be faster as local accesses	Remote access is handled across the network, so slower compared to caching.
2.	It reduces the network traffic, server load and improves scalability as server is contacted occasionally.	It increases network traffic, server load and degrades performance.
3.	In case of caching, for transmission of big chunks of data, overall network overload required is at lower side compared to remote.	In case of remote service, transmission of series of responses to specific requests involve higher network overhead as compared to caching.

Sr. No	Caching	Remote Service
4.	<i>I Caching is better if J case of infrequent I writes but if writes I are frequent then I mechanism used to overcome J consistency I problem incur large I overhead in terms I of performance. network traffic, and server load.</i>	In remote service, there is always communication between client and server to have a master copy consistent with client's cached copy.
5.	Caching is better option for machines with disk or large main memory.	If machines are diskless and with small main memory then remote service should be carried out.
6.	<i>■ In case of caching, the lower-level inter-machine interface is I different from the I "ppcr-Icvc1 user I interface.</i>	The remote service concept is just an extension of the local file-system interface across the network.

Syllabus Topic : Stateful Versus Stateless Service

6JO Stateful Versus Stateless Service

- Client access the remote files from server. If server keeps track about each file being accessed by each client then the service is called stateful. If server simply provides requested blocks of data to client and does not keep track about how client makes use of them then service is called stateless.

Stateful File Service

- First client must carry out open() operation on file prior to accessing it. Server stores the access information about file from disk and store in main memory.
- Server then gives unique connection identifier to client and opens the file for client. This identifier is used by client to access the file throughout the session.
- On closing of the file, or by garbage collection mechanism, the server should free the main memory space used by client when it is no longer active. The information maintained by server regarding client is used for fault tolerance. AFS uses stateful approach

<r Stateless File Service

- in this case server does not keep any information in main memory about opened files. Here, each request identifies the file. There is no need to open and close the file with operations open() and close().
- Each file operation in this case is not a part of session but it is on its own separately. Closing the file at last is also a remote operation, NFS is example of stateless file service approach.
- Stateful service gives more performance. As connection identifier is used to access information maintained in main memory, disk accesses are reduced. Moreover, as server knows that file is open for sequential access, then read ahead next block improves the performance. In stateful case, if server crashes then recovery of volatile state information is complex. A dialog with client is used by recovery protocol. Server also needs to know about crash of client to free the memory space allocated to it. All the operation underway during crash should be aborted.
- In case of stateless service, the effect of server crash and recovery is invisible. Client keeps retransmitting the requests if server crash and no response from it.

Syllabus Topic : File Replication

6J1 File Replication

- Replicating the file on many machines improves availability. If nearest replica is used then service time also reduces. The replica of the same file should be kept on failure independent machines so that availability of one replica is not affected by availability of remaining other replica. Replication of files should be hidden from users.
- It is the responsibility of naming scheme to map a replicated file name to a particular replica. Higher levels should remain invisible from existence of replicas.
- At lower-level different lower-level names are used for different replicas. Replication control such as determination of the degree of replication and placement of replicas should be provided to higher levels. As one replica updates then from user's point of view, the changes should be reflected in other copies as well.

6.12 Distributed Synchronization

centralized synchronization is ...
I. J. to distributed environment also ...

6.12.1 Mutual Exclusion

for different algorithms for mutual exclusion in distributed system.

U T recent algorithm, to achieve mutual exclusion in distributed environment. There are several processes in system numbered from 1 to n, and each process has its own processor.

6.12.1(A) Centralized Approach

Q. Explain centralized algorithm for mutual exclusion in distributed system.

In this algorithm, one of the processes acts as coordinator. Each process wants to enter in critical section sends request message to coordinator. If process then receives reply message from coordinator then proceed to enter in critical section. Once process finish execution in critical section, it sends release message to coordinator and continue its execution.

After receiving the request, first coordinator checks whether critical section is empty or not. If empty then coordinator sends reply message to requesting process and process enters in critical section. If other process is already in critical section then request of requesting process is queued.

When coordinator receives release message from process exiting from critical section, one of the processes from queue is chosen to enter in critical section based on some scheduling algorithm. If FIFO strategy is used then no starvation can occur. Three messages request, reply and release are used in this algorithm.

6.12.1(B) Fully Distributed Approach

Q. Explain distributed algorithm for mutual exclusion in distributed system.

In this algorithm all processes are involved in decision to allow the requesting process to enter in critical section. Process P_i which wants to send message to other processes in group sends request (P_i, TS) .

(n - 1) processes, including 8 itself, in this request message TS is a time stamp. Seating the time at receiving this request message.

If process P_i receives request message from other process, it sends reply message. When process enters critical section, then request is reply to request. If requesting process gets reply from all processes, it can enter in critical section.

Whether reply is sent depends on following factors.

- If process receiving request message is in critical section then it defers the reply to requesting process.
- If process receiving request message do not want to enter in critical section then it immediately sends reply to requesting process.
- If process receiving request message also wants to enter in critical section then it compares its time stamp with time stamp of requesting process. If its time stamp is greater than time stamp of incoming request then it immediately sends request message. Otherwise the reply is postponed.

In this algorithm, mutual exclusion and no deadlock is guaranteed. There is no starvation as time stamp is used to ensure FCFS strategy to serve. Exact $2(n-1)$ number of messages are required to exchange to enter in critical section which is less compared to if processes act independently and concurrently.

Following are the limitations of this algorithm.

- Processes should know names of all other processes in group. There is problem when new process joins the group. In this case, processes must receive names of all other processes in group. Suppose, request and reply messages are in transit and new process joins the group, then problem definitely will arise.

If anyone process fails, then requesting process will not receive reply message from it. The monitoring of all processes slate is required and notification regarding fail process need to be sent to other processes in group. This is required so that P_i sends $shOUW$ message to failed process.



- 3. It is necessary to pause (he processes frequently that have not taken entry in their critical section. This is required to guarantee other processes that they want to enter the critical section.)
- And hence, this algorithm is better suited to small and stable sets of cooperating processes.

6.12.1(C) Token-Passing Approach

- In this algorithm, a token which is a special type of message circulated in the system. Process acquiring token can enter in the critical section which ensures a single process in critical section at a time. In this algorithm, processes in the system are assumed as logically organized in a ring structure. This logical arrangement can be based on some parameter, for example, integer value of IP address of the machine on which process is running.
- Token circulates around the ring and process which want to enter in critical section holds the token and enter in critical section. After the process exits its critical section, the token is circulated again in ring. In this scheme, process exiting from critical section passes the token to its neighbour. In this algorithm, unidirectional link guarantees the freedom from starvation.
- One message per entry is required to enter in critical if all the processes want to enter in critical section, if no one want to enter in critical section then infinite number of messages may require. Limitations of this algorithm are:
 1. Token may be lost. In this case election needs to be called to generate new token.
 2. Failure of any process requires establishing a new ring.

Syllabus Topic : Concurrency Control

6.13 Concurrency Control

- This section explains concurrency control schemes for distributed environment. The management of the transactions or subtransactions accessing database on the local site is carried out by transaction manager of distributed database system.
- Transaction can be local or it can be part of the transaction running at other sites. The concurrency

control is the responsibility of transaction manager. It also maintains log for future recovery purpose.

6.13.1 Locking Protocols

- Two-phase locking protocol also works for distributed environment. In this case implementation of lock manager differ compare to single system environment

6.13.1(A) Non replicated Scheme

- This scheme is used when data is not replicated in system. In this scheme, there is a local lock manager on each site to manage the lock and unlock requests for those data items stored in that site. In order to lock the data item X on site Si, process simply send message to lock manager on that site.
- If X is locked in incompatible mode then the request is delayed until that request can be granted. Once it has been decided that the lock request can be approved, the lock manager sends a message back to the initiator showing that the lock request has been granted.
- The implementation of this scheme is easy. It needs only two message transfers for handling lock requests and one message transfer for handling unlock requests. But, deadlock handling is more complex.

6.13.1(B) Single-Coordinator Approach

- This scheme is used when data is replicated. In this scheme, single lock manager resides on particular selected site Si. For locking and unlocking the data item, all requests are sent to this site. Any transaction wanting to lock the data item, sends lock request to Si and lock manager decides whether to grant lock without delay or not.
- If lock request is granted, then lock manager instantly sends message to the site from where request was sent.
- If request is not granted, then it is delayed till it is granted and again message is sent to the requesting site.
- The transaction can read the data item from any one of the sites at which a replica of the data item resides.
- In case of write operation, all the sites having replica should take part in the writing.

Advantages

- Implementation is simple as only two messages are required to handle lock request and one message to handle unlock request. Algorithms used to handle

deadlock on site 8th system can be used " - a lock unlock request are from same's site.

6-17

pM****

single manager reside on "particular",
site si, site Si may be come "processing" due to
processing of ast e request, then "ck due ?
> » * done or processing must stop " fs recovery

another solution to avoid bottle neck is to mnk...
another approach. Lock manager run on several site
and each handles lock and unlock requests for subset of
data items. Since several sites are involved, deadlock
handling is complex.

6.13 1(C) Majority Protocol

Describe majority protocol in distributed system.

This approach is modification of nonreplicated
scheme. In "heme lock manager is pra, c? data
file. The locks for all data or replica of data TM sX

one site has one lock manager running on that site.

If data item is replicated on n number of sites
transaction want to lock this data item. Then it is
necessary to send the lock request to more than one-
half of the n sites. Lock manager decides whether to
grant the lock or should response be delayed. The
transaction does not operate on replicated data item
until it has successfully obtained a lock on a majority of
the replicas of data item.

Advantages

- Avoids the drawback of central control approach.

Disadvantages

Implementation is complex. It requires $2(n/2 + 1)$
messages for handling lock requests and $(n/2 + 1)$
messages for handling unlock requests. As many sites
are involved, different types of implementation is
required for handling the deadlock.

6.13.1(D) Biased Protocol

The working of biased protocol is close to the majority
Protocol. In biased protocol, requests for shared locks
get more favorable action than requests for exclusive
locks.

Each site contains a lock manager at that site.
the locks for all the data items stored at that site.
Shared and exclusive locks are managed to
different manner.

Shared locks

Distributed Systems

To acquire the lock

onsite ht d " X" x , an5Mion s 'Wiy
E , «uv. X " " P,iCa of that data item X. ~

When transaction wants to lock data item X. u
h o ingare Plica of X lock " ager at each site

If * e *que st b -- priest is delayed tv eRntCd then response to that
on read operation ■ S SC eme , m POses less overhead
advantage Eive - mpared to majority protocol. This
there compared t reSU ,tS When more reads are
is incurred a Writts , extra overhead on write
handling com i s many sites involved, deadlock
EcomplexiUes exists.

6.13.1 (E) Primary Copy

In this approach for each data item X there is a primary
copy in which is kept on exactly on one site. For this
data item X, this is called as primary site of X. When
transaction needs to lock data item X, it should request
lock at primary site of X. If the request is not granted
then response to that request is delayed.

The concurrency control for replicated and
nonreplicated data is handled in the same way.
Implementation is straightforward. If primary site of
data item X fails then although replica is available, data
item X is accessible.

6.13.2 Timestamping

- A unique timestamp is given to each transaction based
on which serialization order is decided.

6.13.2(A) Generation of Unique Timestamps

- Centralized and distributed methods are used to
generate the unique timestamp. In centralized approach,
single site distributes the time stamps which are
generated using logical counter or site's own clock.

- In the distributed method, every site generates a local
timestamp which is unique. For this, either a logical
counter or the local clock is used. The global unique
timestamp is generated by concatenation of the
local unique timestamp with the site identifier. While
the site identifier is kept in the least

"X" is site " that the global timestamps

Operating System (MU - Sem 4 - IT)

- generated in one site are not greater than those generated in another site.
- Any site should not generate local timestamps at a faster rate compared to other sites. Otherwise, faster timestamps will be generated than slower ones. A mechanism should be implemented to guarantee that timestamps will be generated in a fair manner across the system.
- If system clocks are used to generate local timestamps then fair generation is possible if all clocks tick at the same speed. Any faster or slower clock may affect fair generation of timestamps.

6.13.2(B) Timestamp-Ordering Scheme

- In centralized approach, there can be cascading rollbacks if no mechanism is applied to stop a transaction from reading a data item value that is not until now committed. The basic timestamp scheme experiences through the unwanted property that conflicts between transactions are resolved through rollbacks, rather than through wait.
- This problem can be resolved by buffering the different read and write operations (delay) until a time when we are guaranteed that these operations can happen without causing aborts.

Syllabus Topic : Deadlock Handling

6.14 Deadlock Handling

- Deadlock prevention, deadlock avoidance and deadlock detection algorithms for single system can be extended to distributed system.

6.14.1 Deadlock Prevention and Avoidance

Q. Explain deadlock prevention and avoidance in distributed system.

- In distributed system, resource-ordering deadlock-prevention technique can be used by simply defining a global ordering among the system resources. A unique number is assigned to each resource in system and process request the resource with unique number i only if it is not holding a resource with a unique number greater than i .
- Banker's algorithm can also be used in distributed system provided one of the processes is selected as process which keeps all the information required to

carry out the banker's algorithm. The global resource-overhead and simple to implement. The banker's algorithm involves large overhead, its implementation is simple. The number of processes is limited, and from the selected process, the banker may be bottleneck. The banker's algorithm does not work well in distributed systems. Due to the presence of a Padlock-prevention scheme that considers a unique priority number for each resource type. A unique priority number is given to each process in order to control the prevention.

These priorities decide whether process P_i should wait for P_j if P_i has highest priority than P_j . This method prevents deadlocks since, for every edge from P_i to P_j in the wait-for graph, P_i has a higher priority than P_j . Thus, a cycle cannot exist. This scheme leads to starvation as always lower priority processes will be rolled back. Following are the two deadlock-prevention schemes based on timestamps.

The wait-die scheme

- This scheme is based on a non-preemptive approach. When process P_i requests a resource presently held by P_j , P_i is permitted to wait only if it has a lesser timestamp than does P_j (that is, P_i is older than P_j). Otherwise, P_i is rolled back (dies).

The wound-wait scheme

- This scheme is based on a preemptive approach and is a complement to the wait-die scheme. When process P_i requests a resource presently held by P_j , P_i is permitted to wait only if it has a larger timestamp than does P_j (that is, P_i is younger than P_j). Otherwise, P_i is rolled back (P_j is wounded by P_i).
- If rolled back process does not get new timestamp, then above both schemes can avoid starvation. The rolled back process will always receive smaller timestamp as time always increases.
- In the wait-die approach, an older process has to wait for a younger one to free its resource. Therefore, the older the process gets, the more it tends to wait. On the other hand, in the wound-wait approach, an older process never waits for a younger process. In both schemes needless rollbacks may take place.

Deadlock Detection

If P' exists, eff $P*$ defi D_S re S Urce al l \ll a tio s slate
 Wait-for source of Ch yPC exist then <We in Wai,
 ... denotes deadlock. Each site maintains it.
 1 S' ait-for graph n \ll leS " w " fM graph
 i X P"**@*, .0Ca S" or Which are running
 in other *** * onal rCques,n * or holding the
 local to that that sue.

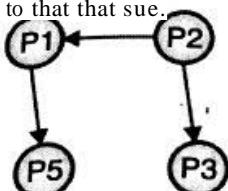


fig. 6-141 : Wait-for graph at site 8,

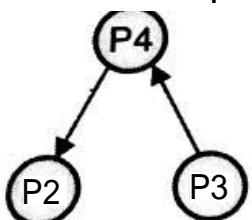


fig. 6.14*2: Wait-for graph at site S2

if process P_i is running on site S_1 and it request for the resource at site S_2 , then P_i sends message to site S_2 and edge P_i to P_j gets added in wait-for graph of site S_2 .

If local wait-for graph contains cycle then deadlock has occurred. If there is no cycle in local wait-for graphs of all sites does not mean deadlock not occurred.

can be a deadlock. The union of the all local wait-for graph is taken on some site to obtain global wait-for graph.

If this global wait-for graph contains cycle then system is in deadlocked state.

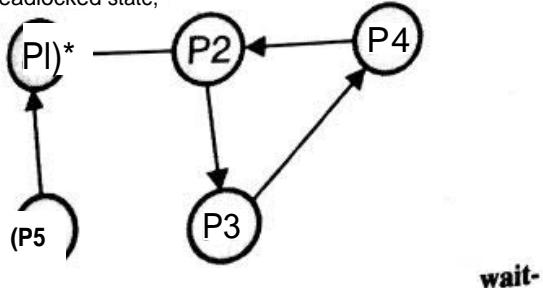


fig. 6.143 : Global wait-for graph after union of wait-for graphs at site S_1 and S_2

Wait-for graphs can be organized in following ways.

Distributed System

6.14.2(A) Centralized Approach

Q. Explain centralized algorithm for deadlock detection in distributed system.

In this approach, there is a deadlock detection coordinator that maintains a global wait-for graph which is union of local wait-for graphs from different sites. The two types of graph are generated. One is real graph which is not illustrated. The constructed graph is an approximation generated by the coordinator during the execution of its algorithm. The reported result by generated graph should be correct. On existence of deadlock, it should report properly and if reported system is certainly in deadlocked state.

The updating of global wait-for graph is carried out if following events take place in system.

- On addition or deletion of edge in local wait-for graphs.
- Periodically, after number of changes.
- On invocation of deadlock detection algorithm.
- On invocation of deadlock detection algorithm, if cycle exist in global wait-for graph then process responsible is rolled back and local sites are also informed by coordinator about the victim. The local sites also roll back the roll out process. On addition or deletion of edge in local wait-for graph, a local site should immediately convey information to coordinator. Coordinator then updates these modifications in global wait-for graph.
- Periodically when local sites sends messages to coordinator about any changes in local wait-for graphs, then false cycle may exists in global wait-for graph.

This happens due to delay incurred in reaching the messages from different local sites to coordinator. Hence, although there is no deadlock, recovery will be initiated. To overcome this problem, messages should be assigned with identifiers or timestamp.

6.14.2(B) Fully Distributed Approach

Q. Explain distributed algorithm for deadlock detection.

In this approach, the sites are involved in deadlock detection.

- In this approach, each site builds a wait-for graph indicating a total graph, tending on the dynamic behavior of the system. The thought is that, if a



deadlock exists, a cycle will appear in at least one of the partial graphs.

In \ll approach, each site has its local graph. In each local wait-for graph, we have $P_i \rightarrow P_j$ for P_i edge P_i to P_j in graph indicates P_i is waiting for P_j to release a resource held in another site being held by any process. In the same way, edge $P_i \rightarrow P_j$ in S_{graph} indicates that a process at another site is waiting to acquire a resource presently being held by P_j in this local site.

These modified wait-for graphs are then used for **detection of deadlock**. If local wait-for graph has cycle

which does not contain P_M then deadlock that involves local processes of that site is occurred. This deadlock is handled locally at that site only.

If local wait-for graph has cycle which contains P_M , then there is possibility of distributed deadlock in which processes from different sites are involved. In this case distributed deadlock detection algorithm is invoked by site whose wait-for graph contains P_M .

6.15 Exam Pack (Review Questions)

Syllabus Topic : Distributed Operating System

- Q. Define distributed system. (Refer section 6.1)
- Q. What are its characteristics? Explain. (Refer section 6.1.1)
- Q. What are objectives behind building the distributed system? (Refer section 6.1.2)
- Q. Explain different types of distributed operating systems. (Refer section 6.2)

Syllabus Topic : Network Based OS

- Q. What are the reasons behind process migration? (Refer section 6.2.2(C))

Syllabus Topic : Network Structure

- Q. Explain LAN type of network. (Refer section 6.3.1)
- Q. Explain WAN type of network. (Refer section 6.3.2)

Syllabus Topic : Network Topology

- Q. Explain different topology for network with its advantages and disadvantages. (Refer section 6.4)

Communication Structure

Syllabus Topic

- Q. What is naming? How names are resolved? Explain. (Refer section 6.5.1)

- Q. Explain different routing strategies. (Refer section 6.5.2)

- Q. What are the different connection strategies used in network communication? (Refer section 6.5.4)

Syllabus Topic : Communication Protocols

- Q. Explain different communication protocols. (Refer section 6.6)

Syllabus Topic : Distributed File Systems

- Q. Explain working of distributed file system (DFS). (Refer section 6.7)

Syllabus Topic : Naming and Transparency

- Q. Explain various naming schemes. (Refer section 6.8.2)

Syllabus Topic : Remote File Access

- Q. What is difference between remote service and caching? (Refer section 6.9.5)

Syllabus Topic : Distributed Synchronization

- Q. Explain centralized algorithm for mutual exclusion in distributed system. (Refer section 6.12.1(A))

- Q. Explain distributed algorithm for mutual exclusion in distributed system. (Refer section 6.12.1(B))

Syllabus Topic : Concurrency Control

- Q. Describe majority protocol in distributed system. (Refer section 6.13.1 (C))

Syllabus Topic : Deadlock Handling

- Q. Explain deadlock prevention and avoidance in distributed system. (Refer section 6.14.1)

- Q. Explain centralized algorithm for deadlock detection in distributed system. (Refer section 6.14.2(B))

- Q. Explain distributed algorithm for deadlock detection in distributed system. (Refer section 6.14.2(C))