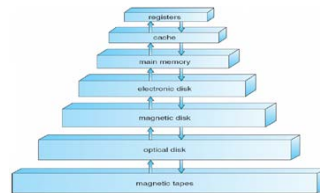


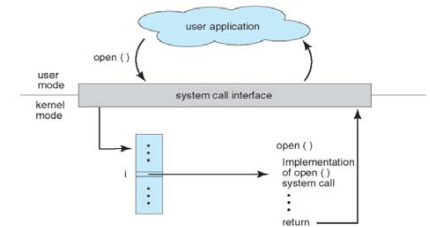
Ch.1 - Introduction

- An OS is a program that acts as an intermediary between a user of a computer and the computer hardware
- Goals: Execute user programs, make the comp. system easy to use, utilize hardware efficiently
- Computer system: Hardware ↔ OS ↔ Applications ↔ Users (↔ = 'uses')
- OS is:
 - Resource allocator: decides between conflicting requests for efficient and fair resource use
 - Control program: controls execution of programs to prevent errors and improper use of computer
- Kernel: the one program running at all times on the computer
- Bootstrap program: loaded at power-up or reboot
 - Stored in ROM or EPROM (known as firmware), Initializes all aspects of system, loads OS kernel and starts execution
- I/O and CPU can execute concurrently
- Device controllers inform CPU that it is finished w/ operation by causing an interrupt
 - Interrupt transfers control to the interrupt service routine generally, through the interrupt vector, which contains the addresses of all the service routines
 - Incoming interrupts are disabled while another interrupt is being processed
 - Trap is a software generated interrupt caused by error or user request
 - OS determines which type of interrupt has occurred by polling or the vectored interrupt system
- System call: request to the operating system to allow user to wait for I/O completion
- Device-status table: contains entry for each I/O device indicating its type, address, and state
 - OS indexes into the I/O device table to determine device status and to modify the table entry to include interrupt
- Storage structure:
 - Main memory – random access, volatile
 - Secondary storage – extension of main memory That provides large non-volatile storage
 - Disk – divided into tracks which are subdivided into sectors. Disk controller determines logical interaction between the device and the computer.
- Caching – copying information into faster storage system
- Multiprocessor Systems: Increased throughput, economy of scale, increased reliability
 - Can be asymmetric or symmetric
 - Clustered systems – Linked multiprocessor systems
- Multiprogramming – Provides efficiency via job scheduling
 - When OS has to wait (ex: for I/O), switches to another job
- Timesharing – CPU switches jobs so frequently that each user can interact with each job while it is running (interactive computing)
- Dual-mode operation allows OS to protect itself and other system components – User mode and kernel mode
 - Some instructions are only executable in kernel mode, these are privileged
- Single-threaded processes have one program counter, multi-threaded processes have one PC per thread
- Protection – mechanism for controlling access of processes or users to resources defined by the OS
- Security – defense of a system against attacks
- User IDs (UID), one per user, and Group IDs, determine which users and groups of users have which privileges



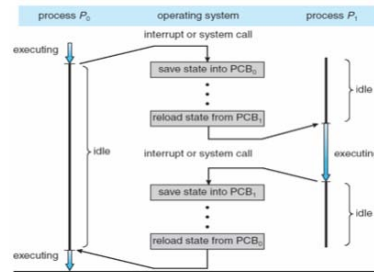
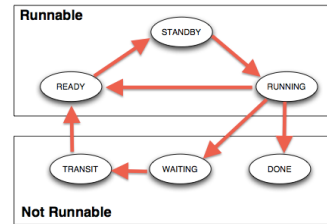
Ch.2 – OS Structures

- User Interface (UI) – Can be Command-Line (CLI) or Graphics User Interface (GUI) or Batch
 - These allow for the user to interact with the system services via system calls (typically written in C/C++)
- Other system services that a helpful to the user include: program execution, I/O operations, file-system manipulation, communications, and error detection
- Services that exist to ensure efficient OS operation are: resource allocation, accounting, protection and security
- Most system calls are accessed by Application Program Interface (API) such as Win32, POSIX, Java
- Usually there is a number associated with each system call
 - System call interface maintains a table indexed according to these numbers
- Parameters may need to be passed to the OS during a system call, may be done by:
 - Passing in registers, address of parameter stored in a block, pushed onto the stack by the program and popped off by the OS
 - Block and stack methods do not limit the number or length of parameters being passed
- Process control system calls include: end, abort, load, execute, create/terminate process, wait, allocate/free memory
- File management system calls include: create/delete file, open/close file, read, write, get/set attributes
- Device management system calls: request/release device, read, write, logically attach/detach devices
- Information maintenance system calls: get/set time, get/set system data, get/set process/file/device attributes
- Communications system calls: create/delete communication connection, send/receive, transfer status information
- OS Layered approach:
 - The operating system is divided into a number of layers (levels), each built on top of lower layers. The bottom layer (layer 0), is the hardware; the highest (layer N) is the user interface
 - With modularity, layers are selected such that each uses functions (operations) and services of only lower-level layers
- Virtual machine: uses layered approach, treats hardware and the OS kernel as though they were all hardware.
 - Host creates the illusion that a process has its own processor and own virtual memory
 - Each guest provided with a 'virtual' copy of the underlying computer
- Application failures can generate core dump file capturing memory of the process
- Operating system failure can generate crash dump file containing kernel memory



Ch.3 – Processes

- Process contains a program counter, stack, and data section.
 - Text section: program code itself
 - Stack: temporary data (function parameters, return addresses, local variables)
 - Data section: global variables
 - Heap: contains memory dynamically allocated during run-time
- Process Control Block (PCB): contains information associated with each process: process state, PC, CPU registers, scheduling information, accounting information, I/O status information
- Types of processes:
 - I/O Bound: spends more time doing I/O than computations, many short CPU bursts
 - CPU Bound: spends more time doing computations, few very long CPU bursts
- When CPU switches to another process, the system must save the state of the old process (to PCB) and load the saved state (from PCB) for the new process via a context switch
 - Time of a context switch is dependent on hardware
- Parent processes create children processes (form a tree)
 - PID allows for process management
 - Parents and children can share all/some/none resources
 - Parents can execute concurrently with children or wait until children terminate
 - fork() system call creates new process
 - exec() system call used after a fork to replace the processes' memory space with a new program
- Cooperating processes need interprocess communication (IPC): shared memory or message passing
- Message passing may be blocking or non-blocking
 - Blocking is considered synchronous
 - Blocking send has the sender block until the message is received
 - Blocking receive has the receiver block until a message is available
 - Non-blocking is considered asynchronous
 - Non-blocking send has the sender send the message and continue
 - Non-blocking receive has the receiver receive a valid message or null



Ch.4 – Threads

- Threads are fundamental unit of CPU utilization that forms the basis of multi-threaded computer systems
- Process creation is heavy-weight while thread creation is light-weight
 - Can simplify code and increase efficiency
- Kernels are generally multi-threaded
- Multi-threading models include: Many-to-One, One-to-One, Many-to-Many
 - Many-to-One: Many user-level threads mapped to single kernel thread
 - One-to-One: Each user-level thread maps to kernel thread
 - Many-to-Many: Many user-level threads mapped to many kernel threads
- Thread library provides programmer with API for creating and managing threads
- Issues include: thread cancellation, signal handling (synchronous/asynchronous), handling thread-specific data, and scheduler activations.
 - Cancellation:
 - Asynchronous cancellation terminates the target thread immediately
 - Deferred cancellation allows the target thread to periodically check if it should be canceled
 - Signal handler processes signals generated by a particular event, delivered to a process, handled
 - Scheduler activations provide upcalls – a communication mechanism from the kernel to the thread library.
 - Allows application to maintain the correct number of kernel threads

Ch.5 – Process Synchronization

- **Race Condition:** several processes access and manipulate the same data concurrently, outcome depends on which order each access takes place.
- Each process has **critical section** of code, where it is manipulating data
 - To solve critical section **problem** each process must ask permission to enter critical section in **entry section**, follow critical section with **exit section** and then execute the **remainder section**
 - Especially difficult to solve this problem in preemptive kernels
- **Peterson's Solution:** solution for two processes

- Two processes share two variables: `int turn` and `Boolean flag[2]`
- **turn:** whose turn it is to enter the critical section
- **flag:** indication of whether or not a process is ready to enter critical section
 - `flag[i] = true` indicates that process P_i is ready
- Algorithm for process P_i :

```
do {
    flag[i] = TRUE;
    turn = j;
    while (flag[j] && turn == j)
        ; // do nothing
    // critical section
    flag[i] = FALSE;
    // remainder section
} while (TRUE);
```

- Modern machines provide atomic hardware instructions: **Atomic** = non-interruptable

- Solution using **Locks**:

```
do {
    acquire lock
    // critical section
    release lock
    // remainder section
} while (TRUE);
```

- Solution using **Test-And-Set**: Shared boolean variable lock, initialized to FALSE

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

```
do {
    while ( TestAndSet (&lock ))
        ; // do nothing
    // critical section
    lock = FALSE;
    // remainder section
} while (TRUE);
```

- Solution using **Swap**: Shared bool variable lock initialized to FALSE; Each process has local bool variable key

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

```
do {
    key = TRUE;
    while ( key == TRUE)
        Swap (&lock, &key);
    // critical section
    lock = FALSE;
    // remainder section
} while (TRUE);
```

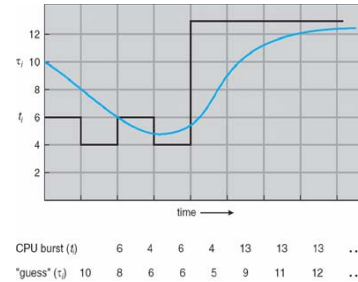
- **Semaphore:** Synchronization tool that does not require busy waiting
 - Standard operations: `wait()` and `signal()` ← these are the only operations that can access semaphore S
 - Can have **counting** (unrestricted range) and **binary** (0 or 1) semaphores
- **Deadlock:** Two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes (most OSes do not prevent or deal with deadlocks)
 - Can cause **starvation** and **priority inversion** (lower priority process holds lock needed by higher-priority process)

Ch.5 – Process Synchronization Continued

- Other synchronization problems include **Bounded-Buffer Problem** and **Readers-Writers Problem**
- **Monitor** is a high-level abstraction that provides a convenient and effective mechanism for process synchronization
 - Only one process may be active within the monitor at a time
 - Can utilize **condition** variables to suspend a resume processes (ex: condition x, y;)
 - `x.wait()` – a process that invokes the operation is suspended until `x.signal()`
 - `x.signal()` – resumes one of processes (if any) that invoked `x.wait()`
 - Can be implemented with semaphores

Ch.6 – CPU Scheduling

- Process execution consists of a cycle of CPU execution and I/O wait
- CPU scheduling decisions take place when a process:
 - Switches from running to waiting (nonpreemptive)
 - Switches from running to ready (preemptive)
 - Switches from waiting to ready (preemptive)
 - Terminates (nonpreemptive)
- The dispatcher module gives control of the CPU to the process selected by the short-term scheduler
 - Dispatch latency- the time it takes for the dispatcher to stop one process and start another
- Scheduling algorithms are chosen based on optimization criteria (ex: throughput, turnaround time, etc.)
 - FCFS, SJF, Shortest-Remaining-Time-First (preemptive SJF), Round Robin, Priority
- Determining length of next CPU burst: Exponential Averaging:
 - t_n = actual length of n^{th} CPU burst
 - τ_{n+1} = predicted value for the next CPU burst
 - $\alpha, 0 \leq \alpha \leq 1$ (commonly α set to $1/2$)
 - Define: $\tau_{n+1} = \alpha * t_n + (1-\alpha)\tau_n$
- Priority Scheduling can result in starvation, which can be solved by aging a process (as time progresses, increase the priority)
- In Round Robin, small time quantum can result in large amounts of context switches
 - Time quantum should be chosen so that 80% of processes have shorter burst times than the time quantum
- Multilevel Queues and Multilevel Feedback Queues have multiple process queues that have different priority levels
 - In the Feedback queue, priority is not fixed → Processes can be promoted and demoted to different queues
 - Feedback queues can have different scheduling algorithms at different levels
- Multiprocessor Scheduling is done in several different ways:
 - Asymmetric multiprocessing: only one processor accesses system data structures → no need to data share
 - Symmetric multiprocessing: each processor is self-scheduling (currently the most common method)
 - Processor affinity: a process running on one processor is more likely to continue to run on the same processor (so that the processor's memory still contains data specific to that specific process)
- Little's Formula can help determine average wait time per process in any scheduling algorithm:
 - $n = \lambda \times W$
 - n = avg queue length; W = avg waiting time in queue; λ = average arrival rate into queue
- Simulations are programmed models of a computer system with variable clocks
 - Used to gather statistics indicating algorithm performance
 - Running simulations is more accurate than queuing models (like Little's Law)
 - Although more accurate, high cost and high risk



Ch.7 – Deadlocks

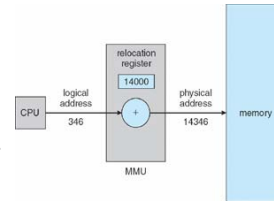
- Deadlock Characteristics: deadlock can occur if these conditions hold simultaneously
 - Mutual Exclusion: only one process at a time can use a resource
 - Hold and Wait: process holding one resource is waiting to acquire resource held by another process
 - No Preemption: a resource can be released only by the process holding it after the process completed its task
 - Circular Wait: set of waiting processes such that P_{n-1} is waiting for resource from P_n , and P_n is waiting for P_0
 - "Dining Philosophers" in deadlock

Ch.8 – Main Memory

- **Cache** sits between main memory and CPU registers
- **Base** and **limit** registers define logical address space usable by a process
- Compiled code addresses **bind** to relocatable addresses
 - Can happen at three different stages
 - **Compile time**: If memory location known a priori, **absolute code** can be generated
 - **Load time**: Must generate **relocatable code** if memory location not known at compile time
 - **Execution time**: Binding delayed until run time if the process can be moved during its execution
- **Memory-Management Unit (MMU)** device that maps virtual to physical address
- Simple scheme uses a **relocation register** which just adds a base value to address
- **Swapping** allows total physical memory space of processes to exceed physical memory
 - Def: process swapped out temporarily to backing store then brought back in for continued execution
- **Backing store**: fast disk large enough to accommodate copies of all memory images
- **Roll out, roll in**: swapping variant for priority-based scheduling.
 - Lower priority process swapped out so that higher priority process can be loaded
- Solutions to **Dynamic Storage-Allocation Problem**:
 - **First-fit**: allocate the first hole that is big enough
 - **Best-fit**: allocate the smallest hole that is big enough (must search entire list) → smallest leftover hole
 - **Worst-fit**: allocate the largest hole (search entire list) → largest leftover hole
- **External Fragmentation**: total memory space exists to satisfy request, but is not contiguous
 - Reduced by **compaction**: relocate free memory to be together in one block
 - Only possible if relocation is dynamic
- **Internal Fragmentation**: allocated memory may be slightly larger than requested memory
- Physical memory divided into fixed-sized frames: size is power of 2, between 512 bytes and 16 MB
- Logical memory divided into same sized blocks: **pages**
- **Page table** used to translate logical to physical addresses
 - **Page number (p)**: used as an index into a page table
 - **Page offset (d)**: combined with base address to define the physical memory address
- **Free-frame list** is maintained to keep track of which frames can be allocated

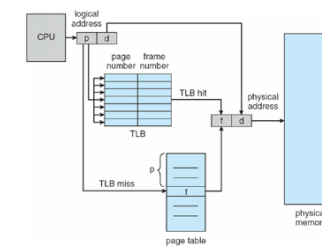


$m - n$ n
For given logical address space 2^m and page size 2^n

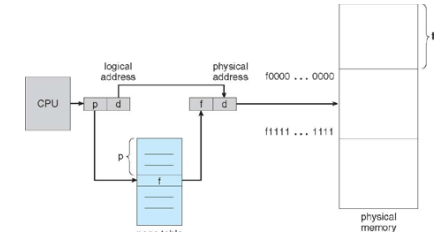


Ch.8 – Main Memory Continued

- **Transition Look-aside Buffer (TLB)** is a CPU cache that memory management hardware uses to improve virtual address translation speed
 - Typically small – 64 to 1024 entries
 - On TLB miss, value loaded to TLB for faster access next time
 - TLB is associative – searched in parallel

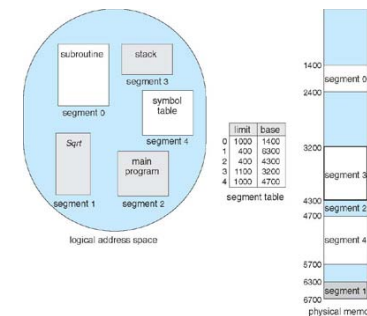


Paging with TLB

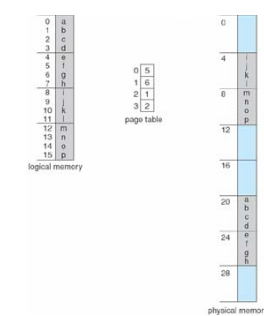


Paging without TLB

- **Effective Access Time**: $EAT = (1 + \epsilon) a + (2 + \epsilon)(1 - \alpha)$
 - ϵ = time unit, α = hit ratio
- **Valid** and **invalid** bits can be used to protect memory
 - “Valid” if the associated page is in the process’ logical address space, so it is a legal page
- Can have multilevel page tables (paged page tables)
- **Hashed Page Tables**: virtual page number hashed into page table
 - Page table has chain of elements hashing to the same location
 - Each element has (1) virtual page number, (2) value of mapped page frame, (3) a pointer to the next element
 - Search through the chain for virtual page number
- **Segment table** – maps two-dimensional physical addresses
 - Entries protected with valid bits and r/w/x privileges



Segmentation example



Page table example

Ch.9 – Virtual Memory

- Virtual memory: separation of user logical memory and physical memory
 - Only part of program needs to be in memory for execution → logical address space > physical address space
 - Allows address spaces to be shared by multiple processes → less swapping
 - Allows pages to be shared during fork(), speeding process creation
- Page fault results from the first time there is a reference to a specific page → traps the OS
 - Must decide to abort if the reference is invalid, or if the desired page is just not in memory yet
 - If the latter: get empty frame, swap page into frame, reset tables to indicate page now in memory, set validation bit, restart instruction that caused the page fault
 - If an instruction accesses multiple pages near each other → less “pain” because of locality of reference
- Demand Paging only brings a page into memory when it is needed → less I/O and memory needed
 - Lazy swapper – never swaps a page into memory unless page will be needed
 - Could result in a lot of page-faults
 - Performance: $EAT = [(1-p)*\text{memory access} + p*(\text{page fault overhead} + \text{swap page out} + \text{swap page in} + \text{restart overhead})]$; where Page Fault Rate $0 \leq p \leq 1$
 - if $p = 0$, no page faults; if $p = 1$, every reference is a fault
 - Can optimize demand paging by loading entire process image to swap space at process load time
- Pure Demand Paging: process starts with no pages in memory
- Copy-on-Write (COW) allows both parent and child processes to initially share the same pages in memory
 - If either process modifies a shared page, only then is the page copied
- Modify (dirty) bit can be used to reduce overhead of page transfers → only modified pages written to disk
- When a page is replaced, write to disk if it has been marked dirty and swap in desired page
- Pages can be replaced using different algorithms: FIFO, LRU (below)
 - Stack can be used to record the most recent page references (LRU is a “stack” algorithm)

reference string

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

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

page frames

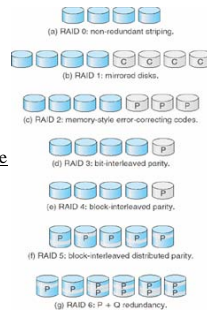
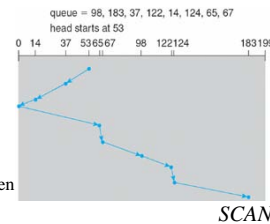
- Second chance algorithm uses a reference bit
 - If 1, decrement and leave in memory
 - If 0, replace next page
- Fixed page allocation: Proportional allocation – Allocate according to size of process
 - s_i = size of process P_i , $S = \sum s_i$, m = total number of frames, a_i – allocation for P_i
 - $a_i = (s_i/S)*m$
- Global replacement: process selects a replacement frame from set of all frames
 - One process can take frame from another
 - Process execution time can vary greatly
 - Greater throughput
- Local replacement: each process selects from only its own set of allocated frames
 - More consistent performance
 - Possible under-utilization of memory
- Page-fault rate is very high if a process does not have “enough” pages
 - Thrashing: a process is busy swapping pages in and out → minimal work is actually being performed
- Memory-mapped file I/O allows file I/O to be treated as routine memory access by mapping a disk block to a page

in memory

- I/O Interlock: Pages must sometimes be locked into memory

Ch.10 – Mass-Storage Systems

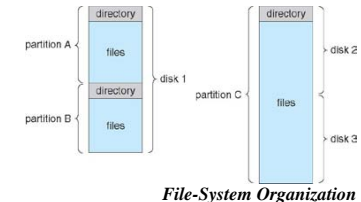
- Magnetic disks provide bulk of secondary storage – rotate at 60 to 250 times per second
 - Transfer rate: rate at which data flows between drive and computer
 - Positioning time (random-access time) is time to move disk arm to desired cylinder (seek time) and time for desired sector to rotate under the disk head (rotational latency)
 - Head crash: disk head making contact with disk surface
- Drive attached to computer's I/O bus – EIDE, ATA, SATA, USB, etc.
 - Host controller uses bus to talk to disk controller
- Access latency = Average access time = average seek time + average latency (fast ~5ms, slow ~14.5ms)
- Average I/O time = avg. access time + (amount to transfer / transfer rate) + controller overhead
 - Ex: to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a .1ms controller overhead = 5ms + 4.17ms + 4KB / 1Gb/sec + 0.1ms = 9.27ms + .12ms = 9.39ms
- Disk drives addressed as 1-dimensional arrays of logical blocks
 - 1-dimensional array is mapped into the sectors of the disk sequentially
- Host-attached storage accessed through I/O ports talking to I/O buses
 - Storage area network (SAN): many hosts attach to many storage units, common in large storage environments
 - Storage made available via LUN masking from specific arrays to specific servers
- Network attached storage (NAS): storage made available over a network rather than local connection
- In disk scheduling, want to minimize seek time; Seek time is proportional to seek distance
- Bandwidth is (total number of bytes transferred) / (total time between first request and completion of last transfer)
- Sources of disk I/O requests: OS, system processes, user processes
 - OS maintains queue of requests, per disk or device
- Several algorithms exist to schedule the servicing of disk I/O requests
 - FCFS, SSTF (shortest seek time first), SCAN, CSCAN, LOOK, CLOOK
 - SCAN/elevator: arm starts at one end and moves towards other end servicing requests as it goes, then reverses direction
 - CSCAN: instead of reversing direction, immediately goes back to beginning
 - LOOK/CLOOK: Arm only goes as far as the last request in each directions, then reverses immediately
- Low level/physical formatting: dividing a disk into sectors that the disk controller can read and write – usually 512 bytes of data
- Partition: divide disk into one or more groups of cylinders, each treated as logical disk
- Logical formatting: “making a file system”
- Increase efficiency by grouping blocks into clusters - Disk I/O is performed on blocks
 - Boot block initializes system - bootstrap loader stored in boot block
- Swap-space: virtual memory uses disk space as an extension of main memory
 - Kernel uses swap maps to track swap space use
- RAID: Multiple disk drives provide reliability via redundancy – increases mean time to failure
 - Disk striping uses group of disks as one storage unit
 - Mirroring/shadowing (RAID 1) – keeps duplicate of each disk
 - Striped mirrors (RAID 1+0) or mirrored striped (RAID 0+1) provides high performance/reliability
 - Block interleaved parity (RAID 4, 5, 6) uses much less redundancy
- Solaris ZFS adds checksums of all data and metadata – detect if object is the right one and whether it changed
- Tertiary storage is usually built using removable media – can be WORM or Read-only, handled like fixed disks
- Fixed disk usually more reliable than removable disk or tape drive
- Main memory is much more expensive than disk storage



Ch.11 – File-System Interface

- File – Uniform logical view of information storage (no matter the medium)
 - Mapped onto physical devices (usually nonvolatile)
 - Smallest allotment of nameable storage
 - Types: Data (numeric, character, binary), Program, Free form, Structured
 - Structure decided by OS and/or program/programmer
- Attributes:
 - Name: Only info in human-readable form
 - Identifier: Unique tag, identifies file within the file system
 - Type, Size
 - Location: pointer to file location
 - Time, date, user identification
- File is an abstract data type
- Operations: create, write, read, reposition within file, delete, truncate
- Global table maintained containing process-independent open file information: open-file table
 - Per-process open file table contains pertinent info, plus pointer to entry in global open file table
- Open file locking: mediates access to a file (shared or exclusive)
 - Mandatory – access denied depending on locks held and requested
 - Advisory – process can find status of locks and decide what to do
- File type can indicate internal file structure
- Access Methods: Sequential access, direct access
 - Sequential Access: tape model of a file
 - Direct Access: random access, relative access
- Disk can be subdivided into partitions; disks or partitions can be RAID protected against failure.
 - Can be used raw without a file-system or formatted with a file system
 - Partitions also known as minidisks, slices
- Volume contains file system: also tracks file system's info in device directory or volume table of contents
- File system can be general or special-purpose. Some special purpose FS:
 - tmpfs – temporary file system in volatile memory
 - objfs – virtual file system that gives debuggers access to kernel symbols
 - ctfcs – virtual file system that maintains info to manage which processes start when system boots
 - lofs – loop back file system allows one file system to be accessed in place of another
 - procsfs – virtual file system that presents information on all processes as a file system
- Directory is similar to symbol table – translating file names into their directory entries
 - Should be efficient, convenient to users, logical grouping
 - Tree structured is most popular – allows for grouping
 - Commands for manipulating: remove – rm<file-name> ; make new sub directory - mkdir<dir-name>
- Current directory: default location for activities – can also specify a path to perform activities in
- Acyclic-graph directories adds ability to directly share directories between users
 - Acyclic can be guaranteed by: only allowing shared files, not shared sub directories; garbage collection; mechanism to check whether new links are OK
- File system must be mounted before it can be accessed – kernel data structure keeps track of mount points
- In a file sharing system User IDs and Group IDs help identify a user's permissions
- Client-server allows multiple clients to mount remote file systems from servers – NFS (UNIX), CIFS (Windows)
- Consistency semantics specify how multiple users are to access a shared file simultaneously – similar to synchronization algorithms from Ch.7
 - One way of protection is Controlled Access: when file created, determine r/w/x access for users/groups

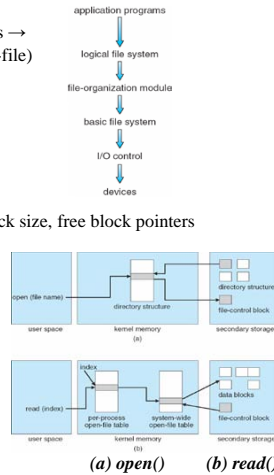
file type	usual extension	function
executable	exe, com, bin or none	ready-to-run machine-language program
object	obj, o	compiled, machine language, not linked
source code	c, cc, java, pas, asm, a	source code in various languages
batch	bat, sh	commands to the command interpreter
text	txt, doc	textual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



File-System Organization

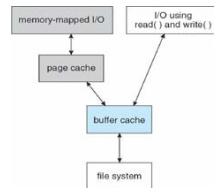
Ch.12 – File System Implementation

- File system resides on secondary storage – disks; file system is organized into layers →
- File control block:** storage structure consisting of information about a file (exist per-file)
- Device driver:** controls the physical device; manage I/O devices
- File organization module:** understands files, logical addresses, and physical blocks
 - Translates logical block number to physical block number
 - Manages free space, disk allocation
- Logical file system:** manages metadata information – maintains file control blocks
- Boot control block:** contains info needed by system to boot OS from volume
- Volume control block:** contains volume details; ex: total # blocks, # free blocks, block size, free block pointers
- Root partition:** contains OS; mounted at boot time
- For all partitions, system is consistency checked at mount time
 - Check metadata for correctness – only allow mount to occur if so
- Virtual file systems provide object-oriented way of implementing file systems
- Directories can be implemented as Linear Lists or Hash Tables
 - Linear list of file names with pointer to data blocks – simple but slow
 - Hash table – linear list with hash data structure – decreased search time
 - Good if entries are fixed size
 - Collisions can occur in hash tables when two file names hash to same location
- Contiguous allocation:** each file occupies set of contiguous blocks
 - Simple, best performance in most cases; problem – finding space for file, external fragmentation
 - Extent based file systems are modified contiguous allocation schemes – extent is allocated for file allocation
- Linked Allocation:** each file is a linked list of blocks – no external fragmentation
 - Locating a block can take many I/Os and disk seeks
- Indexed Allocation:** each file has its own index block(s) of pointers to its data blocks
 - Need index table; can be random access; dynamic access without external fragmentation but has overhead
- Best methods: linked good for sequential, not random; contiguous good for sequential and random
- File system maintains free-space list to track available blocks/clusters
- Bit vector** or **bit map** (n blocks): block number calculation → $(\# \text{bits/word}) * (\# \text{0-value words}) + (\text{offset for } 1^{\text{st}} \text{ bit})$



Example:
 block size = 4KB = 212 bytes
 disk size = 240 bytes (1 terabyte)
 $n = 240/212 = 228$ bits (or 256 MB)
 if clusters of 4 blocks → 64MB of memory

- Space maps (used in ZFS) divide device space into metaslab units and manages metaslabs
 - Each metaslab has associated space map
- Buffer cache** – separate section of main memory for frequently used blocks
- Synchronous** writes sometimes requested by apps or needed by OS – no buffering
- Asynchronous writes are more common, buffer-able, faster
- Free-behind and read-ahead techniques to optimize sequential access
- Page cache** caches pages rather than disk blocks using virtual memory techniques and addresses
 - Memory mapped I/O uses page cache while routine I/O through file system uses buffer (disk) cache
- Unified buffer cache:** uses same page cache to cache both memory-mapped pages and ordinary file system I/O to avoid double caching



Sisteme de Operare

<https://cs.unibuc.ro/~pirofti/>

Mail profesor laborator: mircescu@gmail.com

Punctare

Trebuie promovat laboratorul pentru a ajunge la examen (toate laboratoarele sunt obligatorii).

La laborator se dă și un colocviu.

Opțional: Se va face și un proiect în echipe pentru puncte bonus.

Carte: <https://www.os-book.com/OS9/>

Curs: <https://cs.unibuc.ro/~pirofti/so.html>

Examen: 50p (trebuie luat minim 25p ca să trecem)

- În total 5/6 probleme
- Avem voie cu materiale printate/scrise, nu cu materiale digitale

Model de subiect

- Procese zombie: definiție și exemplu de cod care să producă un proces zombie
- Câte procese și câte fire sunt create?

```
for (i = 0; i < 2; i++) {
    fork()
    pthread_create()
    fork()
}
```

Dacă presupunem că la fork nu se copiază firele de execuție, avem 16 procese și 10 fire.

- Avem o anumită arborescență, care este secvența de cod care o produce?

```
1
/  |  \
2  3  4
pid = fork()
If (pid != 0) {
    pid = fork()
    If (pid != 0) {
        fork()
    }
}
```

- Se dă următoarea listă de procese cu timpii necesari pentru execuție:

Proces	CPU
--------	-----

P0	9
P1	13
P2	8
P3	4
P4	7

Aplicați algoritmul de scheduling Round Robin cu o cuantă de timp $q = 3$

P0 P1 P2 P3 P4 P0 P1 P2 P3 P4 P0 P1 P2
3 6 9 12 15 18 21 24 25 28 31 34

- Se dă următoarea secvență de accesare a paginilor: 1, 2, 3, 4, 5, 6, 1, 2, 3, 4, 5
Avem 3 page frame-uri.

1	1	1	4	4	4	1	1	1	4	4
	2	2	2	5	5	5	2	2	2	5
		3	3	3	6	6	6	3	3	3

Se dă următoarea secvență: 1, 2, 3, 1, 4, 2, 5, 3, 6, 4, 4, 5, 4, 6

1	1	1		1	1	5	5	5	4		4		
	2	2		4	4	4	3	3	3		5		
		3		3	2	2	2	6	6		6		

- Considerați problema filozofilor și soluția propusă mai jos:
do {
 wait(chopstick[i])
 wait(chopstick[(i + 1) % n])
 /* */
 signal(chopstick[i])
 signal(chopstick[(i + 1) % n])
} while(true);
Ce fenomen/problemă apare?

Modificați codul astfel încât filozofi pari iau bețișorul din stânga iar cei impari pe cel din dreapta. Arătați că nu mai apare fenomenul.

Arătați dacă noua soluție satisface cele 3 proprietăți (excluziune mutuală, progres, timp finit de așteptare).

- Fie A o matrice de întregi din $N^{10 \times 10}$, reținută linie cu linie contiguu în memorie.

Avem 3 page frame-uri disponibile.
O pagină de memorie ține 10 întregi.

Prog 1: intră în totalitate într-o pagină P10

```
for (i = 0; i < 10; i++) {  
    for (j = 0; j < 10; j++) {  
        A[i][j] = 0;  
    }  
}
```

Câte pagini are nevoie programul?

10 de la matrice + 1 pentru cod = 11 pagini

Cum arată diagrama Gant?

Codul accesează matricea în ordinea în care e în memorie: $A[0][0]$, ..., $A[0][9]$

Paginile matricei A: 0 1 2 3 4 5 6 7 8 9

Cu 3 frame-uri: tot timpul rețin într-un frame programul, și apoi încarc matricea în celelalte două:

10	10	10	10	10	10	10	10	10	10	10
	0	0	2	2	4	4	6	6	8	8
		1	1	3	3	5	5	7	7	9

```
for (j = 0; j < 10; j++) {  
    for (i = 0; i < 10; i++) {  
        A[i][j] = 0;  
    }  
}
```

Paginile din matrice vor fi accesate secvențial:

1, 2, 3, ..., 7, 8, 9, 1, 2, 3...

Avem de 10 ori mai multe page fault-uri.

Curs opțional

Pentru anul 3: curs opțional de **procesarea semnalelor digitale**, ținut de Irofti.

Introduction



Practice Exercises

- 1.1 What are the three main purposes of an operating system?

Answer:

The three main purposes are:

- To provide an environment for a computer user to execute programs on computer hardware in a convenient and efficient manner.
- To allocate the separate resources of the computer as needed to solve the problem given. The allocation process should be as fair and efficient as possible.
- As a control program it serves two major functions: (1) supervision of the execution of user programs to prevent errors and improper use of the computer, and (2) management of the operation and control of I/O devices.

- 1.2 We have stressed the need for an operating system to make efficient use of the computing hardware. When is it appropriate for the operating system to forsake this principle and to “waste” resources? Why is such a system not really wasteful?

Answer:

Single-user systems should maximize use of the system for the user. A GUI might “waste” CPU cycles, but it optimizes the user’s interaction with the system.

- 1.3 What is the main difficulty that a programmer must overcome in writing an operating system for a real-time environment?

Answer:

The main difficulty is keeping the operating system within the fixed time constraints of a real-time system. If the system does not complete a task in a certain time frame, it may cause a breakdown of the entire system it is running. Therefore when writing an operating system for a real-time system, the writer must be sure that his scheduling schemes don’t allow response time to exceed the time constraint.

- 1.4 Keeping in mind the various definitions of **operating system**, consider whether the operating system should include applications such as Web browsers and mail programs. Argue both that it should and that it should not, and support your answers.

Answer:

An argument in favor of including popular applications with the operating system is that if the application is embedded within the operating system, it is likely to be better able to take advantage of features in the kernel and therefore have performance advantages over an application that runs outside of the kernel. Arguments against embedding applications within the operating system typically dominate however: (1) the applications are applications - and not part of an operating system, (2) any performance benefits of running within the kernel are offset by security vulnerabilities, (3) it leads to a bloated operating system.

- 1.5 How does the distinction between kernel mode and user mode function as a rudimentary form of protection (security) system?

Answer:

The distinction between kernel mode and user mode provides a rudimentary form of protection in the following manner. Certain instructions could be executed only when the CPU is in kernel mode. Similarly, hardware devices could be accessed only when the program is executing in kernel mode. Control over when interrupts could be enabled or disabled is also possible only when the CPU is in kernel mode. Consequently, the CPU has very limited capability when executing in user mode, thereby enforcing protection of critical resources.

- 1.6 Which of the following instructions should be privileged?

- Set value of timer.
- Read the clock.
- Clear memory.
- Issue a trap instruction.
- Turn off interrupts.
- Modify entries in device-status table.
- Switch from user to kernel mode.
- Access I/O device.

Answer:

The following operations need to be privileged: Set value of timer, clear memory, turn off interrupts, modify entries in device-status table, access I/O device. The rest can be performed in user mode.

- 1.7 Some early computers protected the operating system by placing it in a memory partition that could not be modified by either the user job or the operating system itself. Describe two difficulties that you think could arise with such a scheme.

Answer:

The data required by the operating system (passwords, access controls, accounting information, and so on) would have to be stored in or passed through unprotected memory and thus be accessible to unauthorized users.

- 1.8 Some CPUs provide for more than two modes of operation. What are two possible uses of these multiple modes?

Answer:

Although most systems only distinguish between user and kernel modes, some CPUs have supported multiple modes. Multiple modes could be used to provide a finer-grained security policy. For example, rather than distinguishing between just user and kernel mode, you could distinguish between different types of user mode. Perhaps users belonging to the same group could execute each other's code. The machine would go into a specified mode when one of these users was running code. When the machine was in this mode, a member of the group could run code belonging to anyone else in the group.

Another possibility would be to provide different distinctions within kernel code. For example, a specific mode could allow USB device drivers to run. This would mean that USB devices could be serviced without having to switch to kernel mode, thereby essentially allowing USB device drivers to run in a quasi-user/kernel mode.

- 1.9 Timers could be used to compute the current time. Provide a short description of how this could be accomplished.

Answer:

A program could use the following approach to compute the current time using timer interrupts. The program could set a timer for some time in the future and go to sleep. When it is awakened by the interrupt, it could update its local state, which it is using to keep track of the number of interrupts it has received thus far. It could then repeat this process of continually setting timer interrupts and updating its local state when the interrupts are actually raised.

- 1.10 Give two reasons why caches are useful. What problems do they solve? What problems do they cause? If a cache can be made as large as the device for which it is caching (for instance, a cache as large as a disk), why not make it that large and eliminate the device?

Answer:

Caches are useful when two or more components need to exchange data, and the components perform transfers at differing speeds. Caches solve the transfer problem by providing a buffer of intermediate speed between the components. If the fast device finds the data it needs in the cache, it need not wait for the slower device. The data in the cache must be kept consistent with the data in the components. If a component has a data value change, and the datum is also in the cache, the cache must also be updated. This is especially a problem on multiprocessor systems where more than one process may be accessing a datum. A component may be eliminated by an equal-sized cache, but only if: (a) the cache and the component have equivalent state-saving capacity (that is, if the component retains its data when electricity is removed, the cache must

retain data as well), and (b) the cache is affordable, because faster storage tends to be more expensive.

- 1.11 Distinguish between the client-server and peer-to-peer models of distributed systems.

Answer:

The client-server model firmly distinguishes the roles of the client and server. Under this model, the client requests services that are provided by the server. The peer-to-peer model doesn't have such strict roles. In fact, all nodes in the system are considered peers and thus may act as *either* clients or servers—or both. A node may request a service from another peer, or the node may in fact provide such a service to other peers in the system.

For example, let's consider a system of nodes that share cooking recipes. Under the client-server model, all recipes are stored with the server. If a client wishes to access a recipe, it must request the recipe from the specified server. Using the peer-to-peer model, a peer node could ask other peer nodes for the specified recipe. The node (or perhaps nodes) with the requested recipe could provide it to the requesting node. Notice how each peer may act as both a client (it may request recipes) and as a server (it may provide recipes).

Operating System Structures



Practice Exercises

2.1 What is the purpose of system calls?

Answer:

System calls allow user-level processes to request services of the operating system.

2.2 What are the five major activities of an operating system with regard to process management?

Answer:

The five major activities are:

- The creation and deletion of both user and system processes
- The suspension and resumption of processes
- The provision of mechanisms for process synchronization
- The provision of mechanisms for process communication
- The provision of mechanisms for deadlock handling

2.3 What are the three major activities of an operating system with regard to memory management?

Answer:

The three major activities are:

- Keep track of which parts of memory are currently being used and by whom.
- Decide which processes are to be loaded into memory when memory space becomes available.
- Allocate and deallocate memory space as needed.

2.4 What are the three major activities of an operating system with regard to secondary-storage management?

Answer:

The three major activities are:

- Free-space management.
- Storage allocation.
- Disk scheduling.

2.5 What is the purpose of the command interpreter? Why is it usually separate from the kernel?

Answer:

It reads commands from the user or from a file of commands and executes them, usually by turning them into one or more system calls. It is usually not part of the kernel since the command interpreter is subject to changes.

2.6 What system calls have to be executed by a command interpreter or shell in order to start a new process?

Answer:

In Unix systems, a *fork* system call followed by an *exec* system call need to be performed to start a new process. The *fork* call clones the currently executing process, while the *exec* call overlays a new process based on a different executable over the calling process.

2.7 What is the purpose of system programs?

Answer:

System programs can be thought of as bundles of useful system calls. They provide basic functionality to users so that users do not need to write their own programs to solve common problems.

2.8 What is the main advantage of the layered approach to system design? What are the disadvantages of using the layered approach?

Answer:

As in all cases of modular design, designing an operating system in a modular way has several advantages. The system is easier to debug and modify because changes affect only limited sections of the system rather than touching all sections of the operating system. Information is kept only where it is needed and is accessible only within a defined and restricted area, so any bugs affecting that data must be limited to a specific module or layer.

2.9 List five services provided by an operating system, and explain how each creates convenience for users. In which cases would it be impossible for user-level programs to provide these services? Explain your answer.

Answer:

The five services are:

- Program execution.** The operating system loads the contents (or sections) of a file into memory and begins its execution. A user-level program could not be trusted to properly allocate CPU time.
- I/O operations.** Disks, tapes, serial lines, and other devices must be communicated with at a very low level. The user need only specify the device and the operation to perform on it, while the system converts that request into device- or controller-specific commands. User-level programs cannot be trusted to access only devices they

should have access to and to access them only when they are otherwise unused.

- c. **File-system manipulation.** There are many details in file creation, deletion, allocation, and naming that users should not have to perform. Blocks of disk space are used by files and must be tracked. Deleting a file requires removing the name file information and freeing the allocated blocks. Protections must also be checked to assure proper file access. User programs could neither ensure adherence to protection methods nor be trusted to allocate only free blocks and deallocate blocks on file deletion.
- d. **Communications.** Message passing between systems requires messages to be turned into packets of information, sent to the network controller, transmitted across a communications medium, and reassembled by the destination system. Packet ordering and data correction must take place. Again, user programs might not coordinate access to the network device, or they might receive packets destined for other processes.
- e. **Error detection.** Error detection occurs at both the hardware and software levels. At the hardware level, all data transfers must be inspected to ensure that data have not been corrupted in transit. All data on media must be checked to be sure they have not changed since they were written to the media. At the software level, media must be checked for data consistency; for instance, whether the number of allocated and unallocated blocks of storage match the total number on the device. There, errors are frequently process-independent (for instance, the corruption of data on a disk), so there must be a global program (the operating system) that handles all types of errors. Also, by having errors processed by the operating system, processes need not contain code to catch and correct all the errors possible on a system.

- 2.10 Why do some systems store the operating system in firmware, while others store it on disk?

Answer:

For certain devices, such as handheld PDAs and cellular telephones, a disk with a file system may not be available for the device. In this situation, the operating system must be stored in firmware.

- 2.11 How could a system be designed to allow a choice of operating systems from which to boot? What would the bootstrap program need to do?

Answer:

Consider a system that would like to run both Windows XP and three different distributions of Linux (e.g., RedHat, Debian, and Mandrake). Each operating system will be stored on disk. During system boot-up, a special program (which we will call the **boot manager**) will determine which operating system to boot into. This means that rather than initially booting to an operating system, the boot manager will first run during system startup. It is this boot manager that is responsible for determining which system to boot into. Typically boot managers must be stored at

certain locations of the hard disk to be recognized during system startup. Boot managers often provide the user with a selection of systems to boot into; boot managers are also typically designed to boot into a default operating system if no choice is selected by the user.

Processes



Practice Exercises

- 3.1 Using the program shown in Figure 3.30, explain what the output will be at Line A.
Answer: The result is still 5 as the child updates its copy of value. When control returns to the parent, its value remains at 5.
- 3.2 Including the initial parent process, how many processes are created by the program shown in Figure 3.31?
Answer: There are 8 processes created.
- 3.3 Original versions of Apple's mobile iOS operating system provided no means of concurrent processing. Discuss three major complications that concurrent processing adds to an operating system.
Answer: FILL
- 3.4 The Sun UltraSPARC processor has multiple register sets. Describe what happens when a context switch occurs if the new context is already loaded into one of the register sets. What happens if the new context is in memory rather than in a register set and all the register sets are in use?
Answer: The CPU current-register-set pointer is changed to point to the set containing the new context, which takes very little time. If the context is in memory, one of the contexts in a register set must be chosen and be moved to memory, and the new context must be loaded from memory into the set. This process takes a little more time than on systems with one set of registers, depending on how a replacement victim is selected.
- 3.5 When a process creates a new process using the `fork()` operation, which of the following state is shared between the parent process and the child process?
- a. Stack

- b. Heap
c. Shared memory segments

Answer:

Only the shared memory segments are shared between the parent process and the newly forked child process. Copies of the stack and the heap are made for the newly created process.

- 3.6 With respect to the RPC mechanism, consider the “exactly once” semantic. Does the algorithm for implementing this semantic execute correctly even if the ACK message back to the client is lost due to a network problem? Describe the sequence of messages and discuss whether “exactly once” is still preserved.

Answer:

The “exactly once” semantics ensure that a remote procedure will be executed exactly once and only once. The general algorithm for ensuring this combines an acknowledgment (ACK) scheme combined with timestamps (or some other incremental counter that allows the server to distinguish between duplicate messages).

The general strategy is for the client to send the RPC to the server along with a timestamp. The client will also start a timeout clock. The client will then wait for one of two occurrences: (1) it will receive an ACK from the server indicating that the remote procedure was performed, or (2) it will time out. If the client times out, it assumes the server was unable to perform the remote procedure so the client invokes the RPC a second time, sending a later timestamp. The client may not receive the ACK for one of two reasons: (1) the original RPC was never received by the server, or (2) the RPC was correctly received—and performed—by the server but the ACK was lost. In situation (1), the use of ACKs allows the server ultimately to receive and perform the RPC. In situation (2), the server will receive a duplicate RPC and it will use the timestamp to identify it as a duplicate so as not to perform the RPC a second time. It is important to note that the server must send a second ACK back to the client to inform the client the RPC has been performed.

- 3.7 Assume that a distributed system is susceptible to server failure. What mechanisms would be required to guarantee the “exactly once” semantics for execution of RPCs?

Answer:

The server should keep track in stable storage (such as a disk log) information regarding what RPC operations were received, whether they were successfully performed, and the results associated with the operations. When a server crash takes place and a RPC message is received, the server can check whether the RPC had been previously performed and therefore guarantee “exactly once” semantics for the execution of RPCs.

Threads



Practice Exercises

- 4.1 Provide three programming examples in which multithreading provides better performance than a single-threaded solution.

Answer:

- A Web server that services each request in a separate thread.
- A parallelized application such as matrix multiplication where different parts of the matrix may be worked on in parallel.
- An interactive GUI program such as a debugger where a thread is used to monitor user input, another thread represents the running application, and a third thread monitors performance.

- 4.2 What are two differences between user-level threads and kernel-level threads? Under what circumstances is one type better than the other?

Answer:

- User-level threads are unknown by the kernel, whereas the kernel is aware of kernel threads.
- On systems using either M:1 or M:N mapping, user threads are scheduled by the thread library and the kernel schedules kernel threads.
- Kernel threads need not be associated with a process whereas every user thread belongs to a process. Kernel threads are generally more expensive to maintain than user threads as they must be represented with a kernel data structure.

- 4.3 Describe the actions taken by a kernel to context-switch between kernel-level threads.

Answer:

Context switching between kernel threads typically requires saving the value of the CPU registers from the thread being switched out and restoring the CPU registers of the new thread being scheduled.

- 4.4 What resources are used when a thread is created? How do they differ from those used when a process is created?

Answer:

Because a thread is smaller than a process, thread creation typically uses fewer resources than process creation. Creating a process requires allocating a process control block (PCB), a rather large data structure. The PCB includes a memory map, list of open files, and environment variables. Allocating and managing the memory map is typically the most time-consuming activity. Creating either a user or kernel thread involves allocating a small data structure to hold a register set, stack, and priority.

- 4.5 Assume that an operating system maps user-level threads to the kernel using the many-to-many model and that the mapping is done through LWPs. Furthermore, the system allows developers to create real-time threads for use in real-time systems. Is it necessary to bind a real-time thread to an LWP? Explain.

Answer:

Yes. Timing is crucial to real-time applications. If a thread is marked as real-time but is not bound to an LWP, the thread may have to wait to be attached to an LWP before running. Consider if a real-time thread is running (is attached to an LWP) and then proceeds to block (i.e. must perform I/O, has been preempted by a higher-priority real-time thread, is waiting for a mutual exclusion lock, etc.) While the real-time thread is blocked, the LWP it was attached to has been assigned to another thread. When the real-time thread has been scheduled to run again, it must first wait to be attached to an LWP. By binding an LWP to a real-time thread you are ensuring the thread will be able to run with minimal delay once it is scheduled.

Process Synchronization



Practice Exercises

- 5.1 In Section 5.4, we mentioned that disabling interrupts frequently can affect the system's clock. Explain why this can occur and how such effects can be minimized.

Answer:

The system clock is updated at every clock interrupt. If interrupts were disabled—particularly for a long period of time—it is possible the system clock could easily lose the correct time. The system clock is also used for scheduling purposes. For example, the time quantum for a process is expressed as a number of clock ticks. At every clock interrupt, the scheduler determines if the time quantum for the currently running process has expired. If clock interrupts were disabled, the scheduler could not accurately assign time quanta. This effect can be minimized by disabling clock interrupts for only very short periods.

- 5.2 Explain why Windows, Linux, and Solaris implement multiple locking mechanisms. Describe the circumstances under which they use spinlocks, mutex locks, semaphores, adaptive mutex locks, and condition variables. In each case, explain why the mechanism is needed.

Answer:

These operating systems provide different locking mechanisms depending on the application developers' needs. Spinlocks are useful for multiprocessor systems where a thread can run in a busy-loop (for a short period of time) rather than incurring the overhead of being put in a sleep queue. Mutexes are useful for locking resources. Solaris 2 uses adaptive mutexes, meaning that the mutex is implemented with a spin lock on multiprocessor machines. Semaphores and condition variables are more appropriate tools for synchronization when a resource must be held for a long period of time, since spinning is inefficient for a long duration.

- 5.3 What is the meaning of the term **busy waiting**? What other kinds of waiting are there in an operating system? Can busy waiting be avoided altogether? Explain your answer.

Answer:

Busy waiting means that a process is waiting for a condition to be satisfied in a tight loop without relinquishing the processor. Alternatively, a process could wait by relinquishing the processor, and block on a condition and wait to be awakened at some appropriate time in the future. Busy waiting can be avoided but incurs the overhead associated with putting a process to sleep and having to wake it up when the appropriate program state is reached.

- 5.4 Explain why spinlocks are not appropriate for single-processor systems yet are often used in multiprocessor systems.

Answer:

Spinlocks are not appropriate for single-processor systems because the condition that would break a process out of the spinlock can be obtained only by executing a different process. If the process is not relinquishing the processor, other processes do not get the opportunity to set the program condition required for the first process to make progress. In a multiprocessor system, other processes execute on other processors and thereby modify the program state in order to release the first process from the spinlock.

- 5.5 Show that, if the `wait()` and `signal()` semaphore operations are not executed atomically, then mutual exclusion may be violated.

Answer:

A wait operation atomically decrements the value associated with a semaphore. If two wait operations are executed on a semaphore when its value is 1, if the two operations are not performed atomically, then it is possible that both operations might proceed to decrement the semaphore value, thereby violating mutual exclusion.

- 5.6 Illustrate how a binary semaphore can be used to implement mutual exclusion among n processes.

Answer:

The n processes share a semaphore, `mutex`, initialized to 1. Each process P_i is organized as follows:

```
do {  
    wait(mutex);  
  
    /* critical section */  
  
    signal(mutex);  
  
    /* remainder section */  
} while (true);
```

CPU Scheduling

CHAPTER 6

Practice Exercises

- 6.1 A CPU-scheduling algorithm determines an order for the execution of its scheduled processes. Given n processes to be scheduled on one processor, how many different schedules are possible? Give a formula in terms of n .

Answer:

$n!$ (n factorial = $n \times n - 1 \times n - 2 \times \dots \times 2 \times 1$).

- 6.2 Explain the difference between preemptive and nonpreemptive scheduling.

Answer:

Preemptive scheduling allows a process to be interrupted in the midst of its execution, taking the CPU away and allocating it to another process. Nonpreemptive scheduling ensures that a process relinquishes control of the CPU only when it finishes with its current CPU burst.

- 6.3 Suppose that the following processes arrive for execution at the times indicated. Each process will run for the amount of time listed. In answering the questions, use nonpreemptive scheduling, and base all decisions on the information you have at the time the decision must be made.

Process	Arrival Time	Burst Time
P_1	0.0	8
P_2	0.4	4
P_3	1.0	1

- What is the average turnaround time for these processes with the FCFS scheduling algorithm?
- What is the average turnaround time for these processes with the SJF scheduling algorithm?
- The SJF algorithm is supposed to improve performance, but notice that we chose to run process P_1 at time 0 because we did not know

that two shorter processes would arrive soon. Compute what the average turnaround time will be if the CPU is left idle for the first 1 unit and then SJF scheduling is used. Remember that processes P_1 and P_2 are waiting during this idle time, so their waiting time may increase. This algorithm could be known as future-knowledge scheduling.

Answer:

- 10.53
- 9.53
- 6.86

Remember that turnaround time is finishing time minus arrival time, so you have to subtract the arrival times to compute the turnaround times. FCFS is 11 if you forget to subtract arrival time.

- 6.4 What advantage is there in having different time-quantum sizes at different levels of a multilevel queueing system?

Answer:

Processes that need more frequent servicing, for instance, interactive processes such as editors, can be in a queue with a small time quantum. Processes with no need for frequent servicing can be in a queue with a larger quantum, requiring fewer context switches to complete the processing, and thus making more efficient use of the computer.

- 6.5 Many CPU-scheduling algorithms are parameterized. For example, the RR algorithm requires a parameter to indicate the time slice. Multilevel feedback queues require parameters to define the number of queues, the scheduling algorithms for each queue, the criteria used to move processes between queues, and so on.

These algorithms are thus really sets of algorithms (for example, the set of RR algorithms for all time slices, and so on). One set of algorithms may include another (for example, the FCFS algorithm is the RR algorithm with an infinite time quantum). What (if any) relation holds between the following pairs of algorithm sets?

- Priority and SJF
- Multilevel feedback queues and FCFS
- Priority and FCFS
- RR and SJF

Answer:

- The shortest job has the highest priority.
- The lowest level of MLFQ is FCFS.
- FCFS gives the highest priority to the job having been in existence the longest.
- None.

- 6.6 Suppose that a scheduling algorithm (at the level of short-term CPU scheduling) favors those processes that have used the least processor time in the recent past. Why will this algorithm favor I/O-bound programs and yet not permanently starve CPU-bound programs?

Answer:

It will favor the I/O-bound programs because of the relatively short CPU burst request by them; however, the CPU-bound programs will not starve because the I/O-bound programs will relinquish the CPU relatively often to do their I/O.

- 6.7 Distinguish between PCS and SCS scheduling.

Answer:

PCS scheduling is done local to the process. It is how the thread library schedules threads onto available LWPs. SCS scheduling is the situation where the operating system schedules kernel threads. On systems using either many-to-one or many-to-many, the two scheduling models are fundamentally different. On systems using one-to-one, PCS and SCS are the same.

- 6.8 Assume that an operating system maps user-level threads to the kernel using the many-to-many model and that the mapping is done through the use of LWPs. Furthermore, the system allows program developers to create real-time threads. Is it necessary to bind a real-time thread to an LWP?

Answer:

Yes, otherwise a user thread may have to compete for an available LWP prior to being actually scheduled. By binding the user thread to an LWP, there is no latency while waiting for an available LWP; the real-time user thread can be scheduled immediately.

- 6.9 The traditional UNIX scheduler enforces an inverse relationship between priority numbers and priorities: the higher the number, the lower the priority. The scheduler recalculates process priorities once per second using the following function:

$$\text{Priority} = (\text{recent CPU usage} / 2) + \text{base}$$

where $\text{base} = 60$ and *recent CPU usage* refers to a value indicating how often a process has used the CPU since priorities were last recalculated.

Assume that recent CPU usage for process P_1 is 40, for process P_2 is 18, and for process P_3 is 10. What will be the new priorities for these three processes when priorities are recalculated? Based on this information, does the traditional UNIX scheduler raise or lower the relative priority of a CPU-bound process?

Answer:

The priorities assigned to the processes are 80, 69, and 65 respectively. The scheduler lowers the relative priority of CPU-bound processes.

Deadlocks



Practice Exercises

- 7.1 List three examples of deadlocks that are not related to a computer-system environment.

Answer:

- Two cars crossing a single-lane bridge from opposite directions.
- A person going down a ladder while another person is climbing up the ladder.
- Two trains traveling toward each other on the same track.

- 7.2 Suppose that a system is in an unsafe state. Show that it is possible for the processes to complete their execution without entering a deadlock state.

Answer:

An unsafe state may not necessarily lead to deadlock, it just means that we cannot guarantee that deadlock will not occur. Thus, it is possible that a system in an unsafe state may still allow all processes to complete without deadlock occurring. Consider the situation where a system has 12 resources allocated among processes P_0 , P_1 , and P_2 . The resources are allocated according to the following policy:

	Max	Current	Need
P_0	10	5	5
P_1	4	2	2
P_2	9	3	6

Currently there are two resources available. This system is in an unsafe state as process P_1 could complete, thereby freeing a total of four resources. But we cannot guarantee that processes P_0 and P_2 can complete. However, it is possible that a process may release resources before requesting any further. For example, process P_2 could release a resource, thereby increasing the total number of resources to five. This

allows process P_0 to complete, which would free a total of nine resources, thereby allowing process P_2 to complete as well.

- 7.3 Consider the following snapshot of a system:

	<u>Allocation</u>	<u>Max</u>	<u>Available</u>
	<i>A B C D</i>	<i>A B C D</i>	<i>A B C D</i>
P_0	0 0 1 2	0 0 1 2	1 5 2 0
P_1	1 0 0 0	1 7 5 0	
P_2	1 3 5 4	2 3 5 6	
P_3	0 6 3 2	0 6 5 2	
P_4	0 0 1 4	0 6 5 6	

Answer the following questions using the banker's algorithm:

- What is the content of the matrix **Need**?
- Is the system in a safe state?
- If a request from process P_1 arrives for (0,4,2,0), can the request be granted immediately?

Answer:

- The values of **Need** for processes P_0 through P_4 respectively are (0, 0, 0, 0), (0, 7, 5, 0), (1, 0, 0, 2), (0, 0, 2, 0), and (0, 6, 4, 2).
- The system is in a safe state? Yes. With **Available** being equal to (1, 5, 2, 0), either process P_0 or P_3 could run. Once process P_3 runs, it releases its resources, which allow all other existing processes to run.
- The request can be granted immediately? This results in the value of **Available** being (1, 1, 0, 0). One ordering of processes that can finish is P_0 , P_2 , P_3 , P_1 , and P_4 .

- 7.4 A possible method for preventing deadlocks is to have a single, higher-order resource that must be requested before any other resource. For example, if multiple threads attempt to access the synchronization objects $A \cdots E$, deadlock is possible. (Such synchronization objects may include mutexes, semaphores, condition variables, and the like.) We can prevent the deadlock by adding a sixth object F . Whenever a thread wants to acquire the synchronization lock for any object $A \cdots E$, it must first acquire the lock for object F . This solution is known as **containment**: the locks for objects $A \cdots E$ are contained within the lock for object F . Compare this scheme with the circular-wait scheme of Section 7.4.4.

Answer:

This is probably not a good solution because it yields too large a scope. It is better to define a locking policy with as narrow a scope as possible.

- 7.5 Prove that the safety algorithm presented in Section 7.5.3 requires an order of $m \times n^2$ operations.

Answer:

Figure 7.1 provides Java code that implement the safety algorithm of the banker's algorithm (the complete implementation of the banker's algorithm is available with the source code download).

```

for (int i = 0; i < n; i++) {
    // first find a thread that can finish
    for (int j = 0; j < n; j++) {
        if (!finish[j]) {
            boolean temp = true;
            for (int k = 0; k < m; k++) {
                if (need[j][k] > work[k])
                    temp = false;
            }

            if (temp) { // if this thread can finish
                finish[j] = true;
                for (int x = 0; x < m; x++)
                    work[x] += work[j][x];
            }
        }
    }
}

```

Figure 7.1 Banker's algorithm safety algorithm.

As can be seen, the nested outer loops—both of which loop through n times—provide the n^2 performance. Within these outer loops are two sequential inner loops which loop m times. The big-oh of this algorithm is therefore $O(m \times n^2)$.

- 7.6 Consider a computer system that runs 5,000 jobs per month with no deadlock-prevention or deadlock-avoidance scheme. Deadlocks occur about twice per month, and the operator must terminate and rerun about 10 jobs per deadlock. Each job is worth about \$2 (in CPU time), and the jobs terminated tend to be about half-done when they are aborted.

A systems programmer has estimated that a deadlock-avoidance algorithm (like the banker's algorithm) could be installed in the system with an increase in the average execution time per job of about 10 percent. Since the machine currently has 30-percent idle time, all 5,000 jobs per month could still be run, although turnaround time would increase by about 20 percent on average.

- What are the arguments for installing the deadlock-avoidance algorithm?
- What are the arguments against installing the deadlock-avoidance algorithm?

Answer:

An argument for installing deadlock avoidance in the system is that we could ensure deadlock would never occur. In addition, despite the increase in turnaround time, all 5,000 jobs could still run. An argument against installing deadlock avoidance software is that deadlocks occur infrequently and they cost little when they do occur.

- 7.7 Can a system detect that some of its processes are starving? If you answer “yes,” explain how it can. If you answer “no,” explain how the system can deal with the starvation problem.

Answer:

Starvation is a difficult topic to define as it may mean different things for different systems. For the purposes of this question, we will define starvation as the situation whereby a process must wait beyond a reasonable period of time—perhaps indefinitely—before receiving a requested resource. One way of detecting starvation would be to first identify a period of time— T —that is considered unreasonable. When a process requests a resource, a timer is started. If the elapsed time exceeds T , then the process is considered to be starved.

One strategy for dealing with starvation would be to adopt a policy where resources are assigned only to the process that has been waiting the longest. For example, if process P_a has been waiting longer for resource X than process P_b , the request from process P_b would be deferred until process P_a 's request has been satisfied.

Another strategy would be less strict than what was just mentioned. In this scenario, a resource might be granted to a process that has waited less than another process, providing that the other process is not starving. However, if another process is considered to be starving, its request would be satisfied first.

- 7.8 Consider the following resource-allocation policy. Requests for and releases of resources are allowed at any time. If a request for resources cannot be satisfied because the resources are not available, then we check any processes that are blocked waiting for resources. If a blocked process has the desired resources, then these resources are taken away from it and are given to the requesting process. The vector of resources for which the blocked process is waiting is increased to include the resources that were taken away.

For example, consider a system with three resource types and the vector *Available* initialized to (4,2,2). If process P_0 asks for (2,2,1), it gets them. If P_1 asks for (1,0,1), it gets them. Then, if P_0 asks for (0,0,1), it is blocked (resource not available). If P_2 now asks for (2,0,0), it gets the available one (1,0,0) and one that was allocated to P_0 (since P_0 is blocked). P_0 's *Allocation* vector goes down to (1,2,1), and its *Need* vector goes up to (1,0,1).

- Can deadlock occur? If you answer “yes,” give an example. If you answer “no,” specify which necessary condition cannot occur.
- Can indefinite blocking occur? Explain your answer.

Answer:

- Deadlock cannot occur because preemption exists.
- Yes. A process may never acquire all the resources it needs if they are continuously preempted by a series of requests such as those of process C.

- 7.9 Suppose that you have coded the deadlock-avoidance safety algorithm and now have been asked to implement the deadlock-detection algorithm. Can you do so by simply using the safety algorithm code and redefining $\mathbf{Max}[i] = \mathbf{Waiting}[i] + \mathbf{Allocation}[i]$, where $\mathbf{Waiting}[i]$ is a vector specifying the resources for which process i is waiting and $\mathbf{Allocation}[i]$ is as defined in Section 7.5? Explain your answer.

Answer:

Yes. The *Max* vector represents the maximum request a process may make. When calculating the safety algorithm we use the *Need* matrix, which represents $Max - Allocation$. Another way to think of this is $Max = Need + Allocation$. According to the question, the *Waiting* matrix fulfills a role similar to the *Need* matrix, therefore $Max = Waiting + Allocation$.

- 7.10 Is it possible to have a deadlock involving only one single-threaded process? Explain your answer.

Answer:

No. This follows directly from the hold-and-wait condition.

Main Memory



Practice Exercises

- 8.1 Name two differences between logical and physical addresses.

Answer:

A logical address does not refer to an actual existing address; rather, it refers to an abstract address in an abstract address space. Contrast this with a physical address that refers to an actual physical address in memory. A logical address is generated by the CPU and is translated into a physical address by the memory management unit (MMU). Therefore, physical addresses are generated by the MMU.

- 8.2 Consider a system in which a program can be separated into two parts: code and data. The CPU knows whether it wants an instruction (instruction fetch) or data (data fetch or store). Therefore, two base-limit register pairs are provided: one for instructions and one for data. The instruction base-limit register pair is automatically read-only, so programs can be shared among different users. Discuss the advantages and disadvantages of this scheme.

Answer:

The major advantage of this scheme is that it is an effective mechanism for code and data sharing. For example, only one copy of an editor or a compiler needs to be kept in memory, and this code can be shared by all processes needing access to the editor or compiler code. Another advantage is protection of code against erroneous modification. The only disadvantage is that the code and data must be separated, which is usually adhered to in a compiler-generated code.

- 8.3 Why are page sizes always powers of 2?

Answer:

Recall that paging is implemented by breaking up an address into a page and offset number. It is most efficient to break the address into X page bits and Y offset bits, rather than perform arithmetic on the address to calculate the page number and offset. Because each bit position represents a power of 2, splitting an address between bits results in a page size that is a power of 2.

- 8.4 Consider a logical address space of 64 pages of 1024 words each, mapped onto a physical memory of 32 frames.
- How many bits are there in the logical address?
 - How many bits are there in the physical address?

Answer:

- Logical address: 16 bits
- Physical address: 15 bits

- 8.5 What is the effect of allowing two entries in a page table to point to the same page frame in memory? Explain how this effect could be used to decrease the amount of time needed to copy a large amount of memory from one place to another. What effect would updating some byte on the one page have on the other page?

Answer:

By allowing two entries in a page table to point to the same page frame in memory, users can share code and data. If the code is reentrant, much memory space can be saved through the shared use of large programs such as text editors, compilers, and database systems. “Copying” large amounts of memory could be effected by having different page tables point to the same memory location.

However, sharing of nonreentrant code or data means that any user having access to the code can modify it and these modifications would be reflected in the other user’s “copy.”

- 8.6 Describe a mechanism by which one segment could belong to the address space of two different processes.

Answer:

Since segment tables are a collection of base-limit registers, segments can be shared when entries in the segment table of two different jobs point to the same physical location. The two segment tables must have identical base pointers, and the shared segment number must be the same in the two processes.

- 8.7 Sharing segments among processes without requiring that they have the same segment number is possible in a dynamically linked segmentation system.

- Define a system that allows static linking and sharing of segments without requiring that the segment numbers be the same.
- Describe a paging scheme that allows pages to be shared without requiring that the page numbers be the same.

Answer:

Both of these problems reduce to a program being able to reference both its own code and its data without knowing the segment or page number associated with the address. MULTICS solved this problem by associating four registers with each process. One register had the address of the current program segment, another had a base address for the stack, another had a base address for the global data, and so on. The idea is

that all references have to be indirect through a register that maps to the current segment or page number. By changing these registers, the same code can execute for different processes without the same page or segment numbers.

- 8.8 In the IBM/370, memory protection is provided through the use of *keys*. A key is a 4-bit quantity. Each 2K block of memory has a key (the storage key) associated with it. The CPU also has a key (the protection key) associated with it. A store operation is allowed only if both keys are equal, or if either is zero. Which of the following memory-management schemes could be used successfully with this hardware?
- a. Bare machine
 - b. Single-user system
 - c. Multiprogramming with a fixed number of processes
 - d. Multiprogramming with a variable number of processes
 - e. Paging
 - f. Segmentation

Answer:

- a. Protection not necessary, set system key to 0.
- b. Set system key to 0 when in supervisor mode.
- c. Region sizes must be fixed in increments of 2k bytes, allocate key with memory blocks.
- d. Same as above.
- e. Frame sizes must be in increments of 2k bytes, allocate key with pages.
- f. Segment sizes must be in increments of 2k bytes, allocate key with segments.

Virtual Memory



Practice Exercises

- 9.1 Under what circumstances do page faults occur? Describe the actions taken by the operating system when a page fault occurs.

Answer:

A page fault occurs when an access to a page that has not been brought into main memory takes place. The operating system verifies the memory access, aborting the program if it is invalid. If it is valid, a free frame is located and I/O is requested to read the needed page into the free frame. Upon completion of I/O, the process table and page table are updated and the instruction is restarted.

- 9.2 Assume that you have a page-reference string for a process with m frames (initially all empty). The page-reference string has length p ; n distinct page numbers occur in it. Answer these questions for any page-replacement algorithms:

- What is a lower bound on the number of page faults?
- What is an upper bound on the number of page faults?

Answer:

- n
- p

- 9.3 Consider the page table shown in Figure 9.30 for a system with 12-bit virtual and physical addresses and with 256-byte pages. The list of free page frames is D, E, F (that is, D is at the head of the list, E is second, and F is last).

Convert the following virtual addresses to their equivalent physical addresses in hexadecimal. All numbers are given in hexadecimal. (A dash for a page frame indicates that the page is not in memory.)

- 9EF
- 111

- 700
- 0FF

Answer:

- 9EF - 0EF
- 111 - 211
- 700 - D00
- 0FF - EFF

- 9.4 Consider the following page-replacement algorithms. Rank these algorithms on a five-point scale from “bad” to “perfect” according to their page-fault rate. Separate those algorithms that suffer from Belady’s anomaly from those that do not.

- LRU replacement
- FIFO replacement
- Optimal replacement
- Second-chance replacement

Answer:

Rank	Algorithm	Suffer from Belady’s anomaly
1	Optimal	no
2	LRU	no
3	Second-chance	yes
4	FIFO	yes

- 9.5 Discuss the hardware support required to support demand paging.

Answer:

For every memory-access operation, the page table needs to be consulted to check whether the corresponding page is resident or not and whether the program has read or write privileges for accessing the page. These checks have to be performed in hardware. A TLB could serve as a cache and improve the performance of the lookup operation.

- 9.6 Consider the two-dimensional array A:

```
int A[] [] = new int[100][100];
```

where $A[0][0]$ is at location 200 in a paged memory system with pages of size 200. A small process that manipulates the matrix resides in page 0 (locations 0 to 199). Thus, every instruction fetch will be from page 0.

For three page frames, how many page faults are generated by the following array-initialization loops, using LRU replacement and assuming that page frame 1 contains the process and the other two are initially empty?

```

a. for (int j = 0; j < 100; j++)
    for (int i = 0; i < 100; i++)
        A[i][j] = 0;

b. for (int i = 0; i < 100; i++)
    for (int j = 0; j < 100; j++)
        A[i][j] = 0;

```

Answer:

- a. 5,000
- b. 50

9.7 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 would occur for the following replacement algorithms, assuming one, two, three, four, five, six, or seven frames? Remember all frames are initially empty, so your first unique pages will all cost one fault each.

- LRU replacement
- FIFO replacement
- Optimal replacement

Answer:

<u>Number of frames</u>	<u>LRU</u>	<u>FIFO</u>	<u>Optimal</u>
1	20	20	20
2	18	18	15
3	15	16	11
4	10	14	8
5	8	10	7
6	7	10	7
7	7	7	7

9.8 Suppose that you want to use a paging algorithm that requires a reference bit (such as second-chance replacement or working-set model), but the hardware does not provide one. Sketch how you could simulate a reference bit even if one were not provided by the hardware, or explain why it is not possible to do so. If it is possible, calculate what the cost would be.

Answer:

You can use the valid/invalid bit supported in hardware to simulate the reference bit. Initially set the bit to invalid. On first reference a trap to the operating system is generated. The operating system will set a software bit to 1 and reset the valid/invalid bit to valid.

9.9 You have devised a new page-replacement algorithm that you think may be optimal. In some contorted test cases, Belady's anomaly occurs. Is the new algorithm optimal? Explain your answer.

Answer:

No. An optimal algorithm will not suffer from Belady's anomaly because—by definition—an optimal algorithm replaces the page that will not be used for the longest time. Belady's anomaly occurs when a page-replacement algorithm evicts a page that will be needed in the immediate future. An optimal algorithm would not have selected such a page.

9.10 Segmentation is similar to paging but uses variable-sized "pages." Define two segment-replacement algorithms based on FIFO and LRU page-replacement schemes. Remember that since segments are not the same size, the segment that is chosen to be replaced may not be big enough to leave enough consecutive locations for the needed segment. Consider strategies for systems where segments cannot be relocated, and those for systems where they can.

Answer:

- a. **FIFO.** Find the first segment large enough to accommodate the incoming segment. If relocation is not possible and no one segment is large enough, select a combination of segments whose memories are contiguous, which are "closest to the first of the list" and which can accommodate the new segment. If relocation is possible, rearrange the memory so that the first N segments large enough for the incoming segment are contiguous in memory. Add any leftover space to the free-space list in both cases.
- b. **LRU.** Select the segment that has not been used for the longest period of time and that is large enough, adding any leftover space to the free space list. If no one segment is large enough, select a combination of the "oldest" segments that are contiguous in memory (if relocation is not available) and that are large enough. If relocation is available, rearrange the oldest N segments to be contiguous in memory and replace those with the new segment.

9.11 Consider a demand-paged computer system where the degree of multiprogramming is currently fixed at four. The system was recently measured to determine utilization of CPU and the paging disk. The results are one of the following alternatives. For each case, what is happening? Can the degree of multiprogramming be increased to increase the CPU utilization? Is the paging helping?

- a. CPU utilization 13 percent; disk utilization 97 percent
- b. CPU utilization 87 percent; disk utilization 3 percent
- c. CPU utilization 13 percent; disk utilization 3 percent

Answer:

- a. Thrashing is occurring.
- b. CPU utilization is sufficiently high to leave things alone, and increase degree of multiprogramming.

c. Increase the degree of multiprogramming.

- 9.12 We have an operating system for a machine that uses base and limit registers, but we have modified the machine to provide a page table. Can the page tables be set up to simulate base and limit registers? How can they be, or why can they not be?

Answer:

The page table can be set up to simulate base and limit registers provided that the memory is allocated in fixed-size segments. In this way, the base of a segment can be entered into the page table and the valid/invalid bit used to indicate that portion of the segment as resident in the memory. There will be some problem with internal fragmentation.

Mass Storage Structure



Practice Exercises

- 10.1 Is disk scheduling, other than FCFS scheduling, useful in a single-user environment? Explain your answer.
Answer:
In a single-user environment, the I/O queue usually is empty. Requests generally arrive from a single process for one block or for a sequence of consecutive blocks. In these cases, FCFS is an economical method of disk scheduling. But LOOK is nearly as easy to program and will give much better performance when multiple processes are performing concurrent I/O, such as when a Web browser retrieves data in the background while the operating system is paging and another application is active in the foreground.
- 10.2 Explain why SSTF scheduling tends to favor middle cylinders over the innermost and outermost cylinders.
Answer:
The center of the disk is the location having the smallest average distance to all other tracks. Thus the disk head tends to move away from the edges of the disk. Here is another way to think of it. The current location of the head divides the cylinders into two groups. If the head is not in the center of the disk and a new request arrives, the new request is more likely to be in the group that includes the center of the disk; thus, the head is more likely to move in that direction.
- 10.3 Why is rotational latency usually not considered in disk scheduling? How would you modify SSTF, SCAN, and C-SCAN to include latency optimization?
Answer:
Most disks do not export their rotational position information to the host. Even if they did, the time for this information to reach the scheduler would be subject to imprecision and the time consumed by the scheduler is variable, so the rotational position information would become incorrect. Further, the disk requests are usually given in terms

of logical block numbers, and the mapping between logical blocks and physical locations is very complex.

- 10.4 Why is it important to balance file system I/O among the disks and controllers on a system in a multitasking environment?

Answer:

A system can perform only at the speed of its slowest bottleneck. Disks or disk controllers are frequently the bottleneck in modern systems as their individual performance cannot keep up with that of the CPU and system bus. By balancing I/O among disks and controllers, neither an individual disk nor a controller is overwhelmed, so that bottleneck is avoided.

- 10.5 What are the tradeoffs involved in rereading code pages from the file system versus using swap space to store them?

Answer:

If code pages are stored in swap space, they can be transferred more quickly to main memory (because swap space allocation is tuned for faster performance than general file system allocation). Using swap space can require startup time if the pages are copied there at process invocation rather than just being paged out to swap space on demand. Also, more swap space must be allocated if it is used for both code and data pages.

- 10.6 Is there any way to implement truly stable storage? Explain your answer.

Answer:

Truly stable storage would never lose data. The fundamental technique for stable storage is to maintain multiple copies of the data, so that if one copy is destroyed, some other copy is still available for use. But for any scheme, we can imagine a large enough disaster that all copies are destroyed.

- 10.7 It is sometimes said that tape is a sequential-access medium, whereas a magnetic disk is a random-access medium. In fact, the suitability of a storage device for random access depends on the transfer size. The term *streaming transfer rate* denotes the rate for a data transfer that is underway, excluding the effect of access latency. By contrast, the *effective transfer rate* is the ratio of total bytes per total seconds, including overhead time such as access latency.

Suppose that, in a computer, the level-2 cache has an access latency of 8 nanoseconds and a streaming transfer rate of 800 megabytes per second, the main memory has an access latency of 60 nanoseconds and a streaming transfer rate of 80 megabytes per second, the magnetic disk has an access latency of 15 milliseconds and a streaming transfer rate of 5 megabytes per second, and a tape drive has an access latency of 60 seconds and a streaming transfer rate of 2 megabytes per seconds.

- a. Random access causes the effective transfer rate of a device to decrease, because no data are transferred during the access time. For the disk described, what is the effective transfer rate if an

average access is followed by a streaming transfer of (1) 512 bytes, (2) 8 kilobytes, (3) 1 megabyte, and (4) 16 megabytes?

- The utilization of a device is the ratio of effective transfer rate to streaming transfer rate. Calculate the utilization of the disk drive for each of the four transfer sizes given in part a.
- Suppose that a utilization of 25 percent (or higher) is considered acceptable. Using the performance figures given, compute the smallest transfer size for disk that gives acceptable utilization.
- Complete the following sentence: A disk is a random-access device for transfers larger than _____ bytes and is a sequential-access device for smaller transfers.
- Compute the minimum transfer sizes that give acceptable utilization for cache, memory, and tape.
- When is a tape a random-access device, and when is it a sequential-access device?

Answer:

- For 512 bytes, the effective transfer rate is calculated as follows.
 $ETR = \text{transfer size} / \text{transfer time}$.
 If X is transfer size, then transfer time is $((X/STR) + \text{latency})$.
 $\text{Transfer time is } 15\text{ms} + (512\text{B}/5\text{MB per second}) = 15.0097\text{ms}$.
 $\text{Effective transfer rate is therefore } 512\text{B}/15.0097\text{ms} = 33.12 \text{ KB/sec}$.
 $ETR \text{ for } 8\text{KB} = .47\text{MB/sec}$.
 $ETR \text{ for } 1\text{MB} = 4.65\text{MB/sec}$.
 $ETR \text{ for } 16\text{MB} = 4.98\text{MB/sec}$.
- Utilization of the device for 512B = $33.12 \text{ KB/sec} / 5\text{MB/sec} = .0064 = .64$
 For 8KB = 9.4%.
 For 1MB = 93%.
 For 16MB = 99.6%.
- Calculate $.25 = ETR/STR$, solving for transfer size X .
 $STR = 5\text{MB}$, so $1.25\text{MB}/S = ETR$.
 $1.25\text{MB}/S * ((X/5) + .015) = X$.
 $.25X + .01875 = X$.
 $X = .025\text{MB}$.
- A disk is a random-access device for transfers larger than K bytes (where $K > \text{disk block size}$), and is a sequential-access device for smaller transfers.
- Calculate minimum transfer size for acceptable utilization of cache memory:
 $STR = 800\text{MB}$, $ETR = 200$, $\text{latency} = 8 * 10^{-9}$.
 $200 (X\text{MB}/800 + 8 * 10^{-9}) = X\text{MB}$.
 $.25X\text{MB} + 1600 * 10^{-9} = X\text{MB}$.
 $X = 2.24 \text{ bytes}$.
 Calculate for memory:

$STR = 80\text{MB}$, $ETR = 20$, $L = 60 * 10^{-9}$.
 $20 (X\text{MB}/80 + 60 * 10^{-9}) = X\text{MB}$.
 $.25X\text{MB} + 1200 * 10^{-9} = X\text{MB}$.
 $X = 1.68 \text{ bytes}$.
 Calculate for tape:
 $STR = 2\text{MB}$, $ETR = .5$, $L = 60\text{s}$.
 $.5 (X\text{MB}/2 + 60) = X\text{MB}$.
 $.25X\text{MB} + 30 = X\text{MB}$.
 $X = 40\text{MB}$.

- It depends upon how it is being used. Assume we are using the tape to restore a backup. In this instance, the tape acts as a sequential-access device where we are sequentially reading the contents of the tape. As another example, assume we are using the tape to access a variety of records stored on the tape. In this instance, access to the tape is arbitrary and hence considered random.

- 10.8** Could a RAID level 1 organization achieve better performance for read requests than a RAID level 0 organization (with nonredundant striping of data)? If so, how?

Answer:

Yes, a RAID Level 1 organization could achieve better performance for read requests. When a read operation is performed, a RAID Level 1 system can decide which of the two copies of the block should be accessed to satisfy the request. This choice could be based on the current location of the disk head and could therefore result in performance optimizations by choosing a disk head that is closer to the target data.

File-System Interface



Practice Exercises

- 11.1 Some systems automatically delete all user files when a user logs off or a job terminates, unless the user explicitly requests that they be kept; other systems keep all files unless the user explicitly deletes them. Discuss the relative merits of each approach.

Answer:

Deleting all files not specifically saved by the user has the advantage of minimizing the file space needed for each user by not saving unwanted or unnecessary files. Saving all files unless specifically deleted is more secure for the user in that it is not possible to lose files inadvertently by forgetting to save them.

- 11.2 Why do some systems keep track of the type of a file, while others leave it to the user and others simply do not implement multiple file types? Which system is “better?”

Answer:

Some systems allow different file operations based on the type of the file (for instance, an ascii file can be read as a stream while a database file can be read via an index to a block). Other systems leave such interpretation of a file’s data to the process and provide no help in accessing the data. The method that is “better” depends on the needs of the processes on the system, and the demands the users place on the operating system. If a system runs mostly database applications, it may be more efficient for the operating system to implement a database-type file and provide operations, rather than making each program implement the same thing (possibly in different ways). For general-purpose systems it may be better to only implement basic file types to keep the operating system size smaller and allow maximum freedom to the processes on the system.

- 11.3 Similarly, some systems support many types of structures for a file’s data, while others simply support a stream of bytes. What are the advantages and disadvantages of each approach?

Answer:

An advantage of having the system support different file structures is that the support comes from the system; individual applications are not required to provide the support. In addition, if the system provides the support for different file structures, it can implement the support presumably more efficiently than an application.

The disadvantage of having the system provide support for defined file types is that it increases the size of the system. In addition, applications that may require different file types other than what is provided by the system may not be able to run on such systems.

An alternative strategy is for the operating system to define no support for file structures and instead treat all files as a series of bytes. This is the approach taken by UNIX systems. The advantage of this approach is that it simplifies the operating system support for file systems, as the system no longer has to provide the structure for different file types. Furthermore, it allows applications to define file structures, thereby alleviating the situation where a system may not provide a file definition required for a specific application.

- 11.4 Could you simulate a multilevel directory structure with a single-level directory structure in which arbitrarily long names can be used? If your answer is yes, explain how you can do so, and contrast this scheme with the multilevel directory scheme. If your answer is no, explain what prevents your simulation’s success. How would your answer change if file names were limited to seven characters?

Answer:

If arbitrarily long names can be used then it is possible to simulate a multilevel directory structure. This can be done, for example, by using the character “.” to indicate the end of a subdirectory. Thus, for example, the name *jim.java.F1* specifies that *F1* is a file in subdirectory *java* which in turn is in the root directory *jim*.

If file names were limited to seven characters, then the above scheme could not be utilized and thus, in general, the answer is *no*. The next best approach in this situation would be to use a specific file as a symbol table (directory) to map arbitrarily long names (such as *jim.java.F1*) into shorter arbitrary names (such as *XX00743*), which are then used for actual file access.

- 11.5 Explain the purpose of the `open()` and `close()` operations.

Answer:

The purpose of the `open()` and `close()` operations is:

- The `open()` operation informs the system that the named file is about to become active.
- The `close()` operation informs the system that the named file is no longer in active use by the user who issued the close operation.

- 11.6 In some systems, a subdirectory can be read and written by an authorized user, just as ordinary files can be.

- a. Describe the protection problems that could arise.

- b. Suggest a scheme for dealing with each of these protection problems.

Answer:

- a. One piece of information kept in a directory entry is file location. If a user could modify this location, then he could access other files defeating the access-protection scheme.
 - b. Do not allow the user to directly write onto the subdirectory. Rather, provide system operations to do so.
- 11.7 Consider a system that supports 5,000 users. Suppose that you want to allow 4,990 of these users to be able to access one file.
- a. How would you specify this protection scheme in UNIX?
 - b. Can you suggest another protection scheme that can be used more effectively for this purpose than the scheme provided by UNIX?

Answer:

- a. There are two methods for achieving this:
 - i. Create an access control list with the names of all 4990 users.
 - ii. Put these 4990 users in one group and set the group access accordingly. This scheme cannot always be implemented since user groups are restricted by the system.
 - b. The universal access to files applies to all users unless their name appears in the access-control list with different access permission. With this scheme you simply put the names of the remaining ten users in the access control list but with no access privileges allowed.
- 11.8 Researchers have suggested that, instead of having an access list associated with each file (specifying which users can access the file, and how), we should have a *user control list* associated with each user (specifying which files a user can access, and how). Discuss the relative merits of these two schemes.

Answer:

- *File control list.* Since the access control information is concentrated in one single place, it is easier to change access control information and this requires less space.
- *User control list.* This requires less overhead when opening a file.

File-System Implementation



Practice Exercises

- 12.1 Consider a file currently consisting of 100 blocks. Assume that the file-control block (and the index block, in the case of indexed allocation) is already in memory. Calculate how many disk I/O operations are required for contiguous, linked, and indexed (single-level) allocation strategies, if, for one block, the following conditions hold. In the contiguous-allocation case, assume that there is no room to grow at the beginning but there is room to grow at the end. Also assume that the block information to be added is stored in memory.
- The block is added at the beginning.
 - The block is added in the middle.
 - The block is added at the end.
 - The block is removed from the beginning.
 - The block is removed from the middle.
 - The block is removed from the end.

Answer:

The results are:

	Contiguous	Linked	Indexed
a.	201	1	1
b.	101	52	1
c.	1	3	1
d.	198	1	0
e.	98	52	0
f.	0	100	0

- 12.2 What problems could occur if a system allowed a file system to be mounted simultaneously at more than one location?

Answer:

There would be multiple paths to the same file, which could confuse users or encourage mistakes (deleting a file with one path deletes the file in all the other paths).

- 12.3 Why must the bit map for file allocation be kept on mass storage, rather than in main memory?

Answer:

In case of system crash (memory failure) the free-space list would not be lost as it would be if the bit map had been stored in main memory.

- 12.4 Consider a system that supports the strategies of contiguous, linked, and indexed allocation. What criteria should be used in deciding which strategy is best utilized for a particular file?

Answer:

- **Contiguous**—if file is usually accessed sequentially, if file is relatively small.
- **Linked**—if file is large and usually accessed sequentially.
- **Indexed**—if file is large and usually accessed randomly.

- 12.5 One problem with contiguous allocation is that the user must preallocate enough space for each file. If the file grows to be larger than the space allocated for it, special actions must be taken. One solution to this problem is to define a file structure consisting of an initial contiguous area (of a specified size). If this area is filled, the operating system automatically defines an overflow area that is linked to the initial contiguous area. If the overflow area is filled, another overflow area is allocated. Compare this implementation of a file with the standard contiguous and linked implementations.

Answer:

This method requires more overhead than the standard contiguous allocation. It requires less overhead than the standard linked allocation.

- 12.6 How do caches help improve performance? Why do systems not use more or larger caches if they are so useful?

Answer:

Caches allow components of differing speeds to communicate more efficiently by storing data from the slower device, temporarily, in a faster device (the cache). Caches are, almost by definition, more expensive than the device they are caching for, so increasing the number or size of caches would increase system cost.

- 12.7 Why is it advantageous for the user for an operating system to dynamically allocate its internal tables? What are the penalties to the operating system for doing so?

Answer:

Dynamic tables allow more flexibility in system use growth — tables are never exceeded, avoiding artificial use limits. Unfortunately, kernel structures and code are more complicated, so there is more potential for bugs. The use of one resource can take away more system resources (by growing to accommodate the requests) than with static tables.

- 12.8 Explain how the VFS layer allows an operating system to support multiple types of file systems easily.

Answer:

VFS introduces a layer of indirection in the file system implementation. In many ways, it is similar to object-oriented programming techniques. System calls can be made generically (independent of file system type). Each file system type provides its function calls and data structures to the VFS layer. A system call is translated into the proper specific functions for the target file system at the VFS layer. The calling program has no file-system-specific code, and the upper levels of the system call structures likewise are file system-independent. The translation at the VFS layer turns these generic calls into file-system-specific operations.