SATHYABAMA
INSTITUTE OF SCIENCE AND TECHNOLOGY
(DEEMED TO BE UNIVERSITY)
Accredited "A" Grade by NAAC | 12B Status by UGC | Approved by AICTE
www.sathyabama.ac.in

# SCS1501 - OPERATING SYSTEM

# UNIT – 4

# MEMORY MANAGEMENT

**Mrs. C. Kavitha**

ASSISTANT PROFESSOR

DEPARTMENT OF CSE

SCHOOL OF COMPUTING

# UNIT 4
# Syllabus

## UNIT 4    MEMORY MANAGEMENT                                      9 Hrs

Storage Management Strategies - Contiguous Vs. Non-Contiguous Storage Allocation - Fixed & Variable Partition Multiprogramming - Paging - Segmentation - Paging/Segmentation Systems - Page Replacement Strategies - Demand & Anticipatory Paging - File Concept - Access Methods - Directory Structure - File Sharing - Protection File - System Structure - Implementation.

# Memory Management

- Background
- Swapping
- Contiguous Memory Allocation
- Segmentation
- Paging
- Structure of the Page Table
- Example: The Intel 32 and 64-bit Architectures
- Example: ARM Architecture

# Objectives

- To provide a detailed description of various ways of organizing memory hardware

- To discuss various memory-management techniques, including paging and segmentation

- To provide a detailed description of the Intel Pentium, which supports both pure segmentation and segmentation with paging
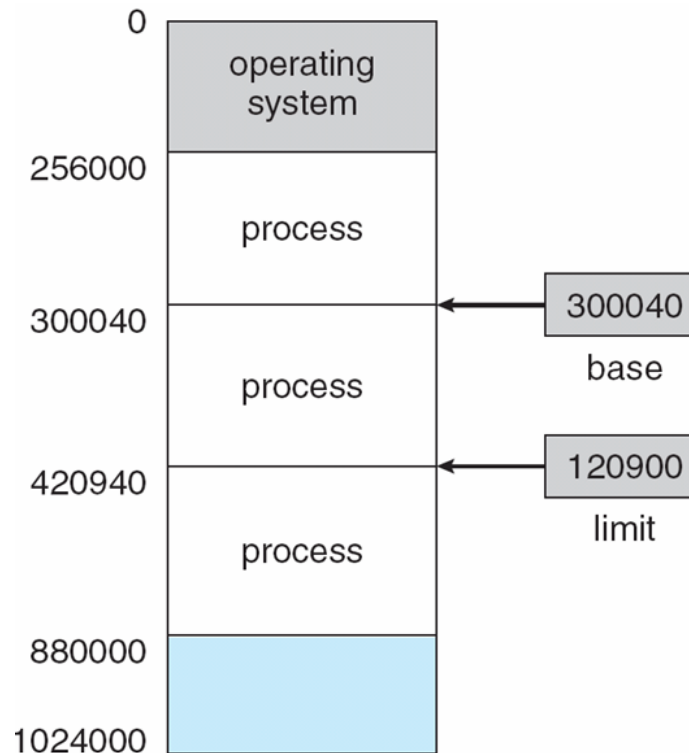
# Background

- Program must be brought (from disk) into memory and placed within a process for it to be run

- Main memory and registers are only storage CPU can access directly

- Memory unit only sees a stream of addresses + read requests, or address + data and write requests

- Register access in one CPU clock (or less)

- Main memory can take many cycles, causing a **stall**

- **Cache** sits between main memory and CPU registers

- Protection of memory required to ensure correct operation
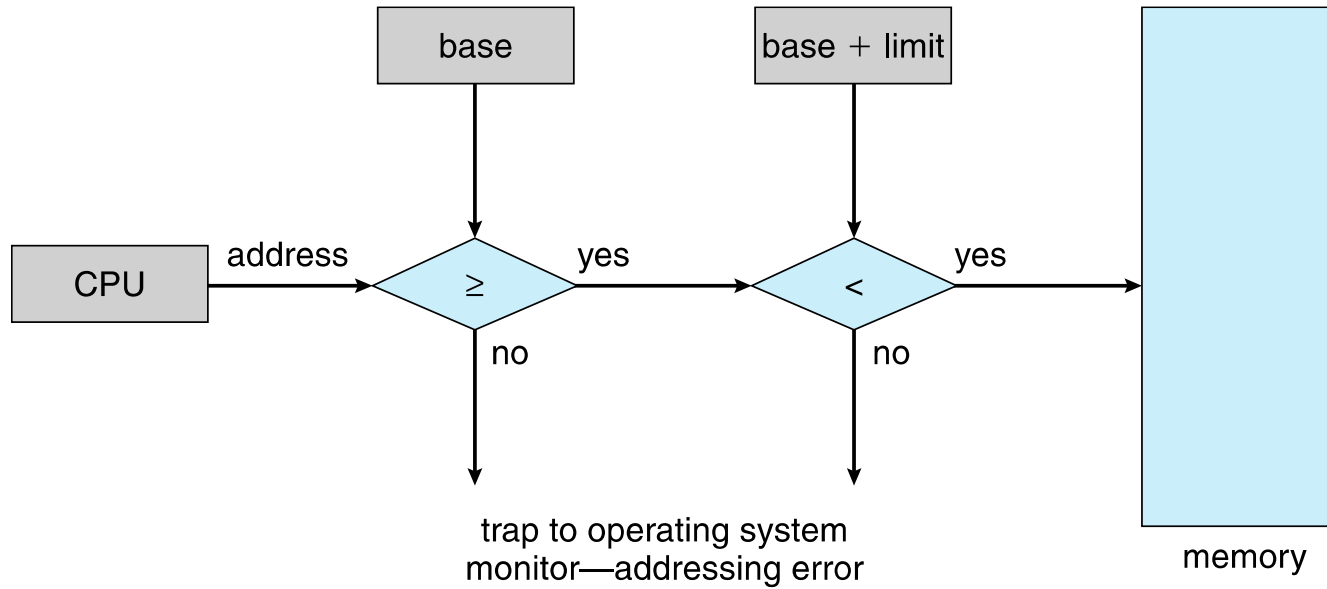
# Base and Limit Registers

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

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

# Hardware Address Protection

# Address Binding

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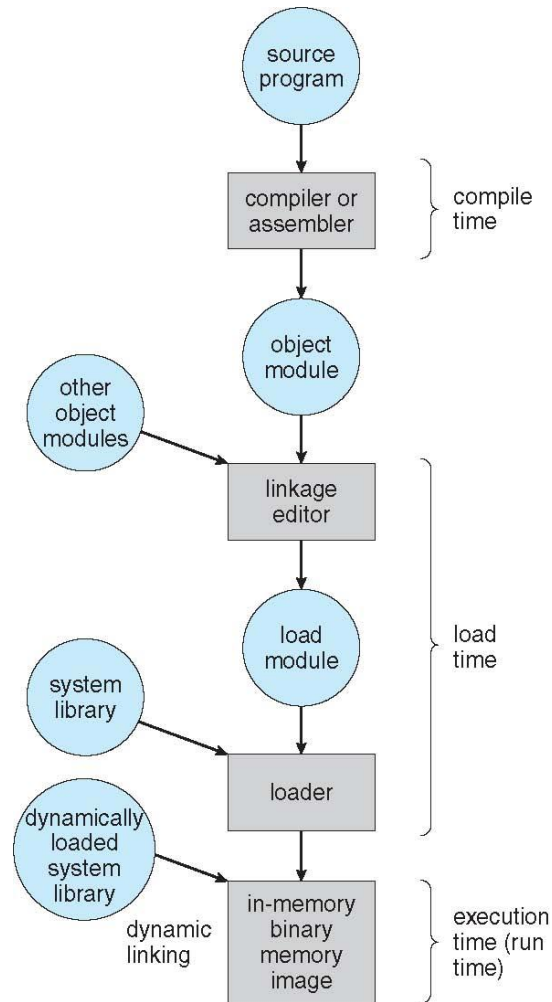    - Each binding maps one address space to another

# Binding of Instructions and Data to Memory

- Address binding of instructions and data to memory addresses can happen at three different stages

  - **Compile time**: If memory location known a priori, **absolute code** can be generated; must recompile code if starting location changes

  - **Load time**: Must generate **relocatable code** if memory location is not known at compile time

  - **Execution time**: Binding delayed until run time if the process can be moved during its execution from one memory segment to another

    - Need hardware support for address maps (e.g., base and limit registers)

# Multistep Processing of a User Program

# Logical vs. Physical Address Space

- The concept of a logical address space that is bound to a separate **physical address space** is central to proper memory management

  - **Logical address** – generated by the CPU; also referred to as **virtual address**

  - **Physical address** – address seen by the memory unit

- Logical and physical addresses are the same in compile-time and load-time address-binding schemes; logical (virtual) and physical addresses differ in execution-time address-binding scheme

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

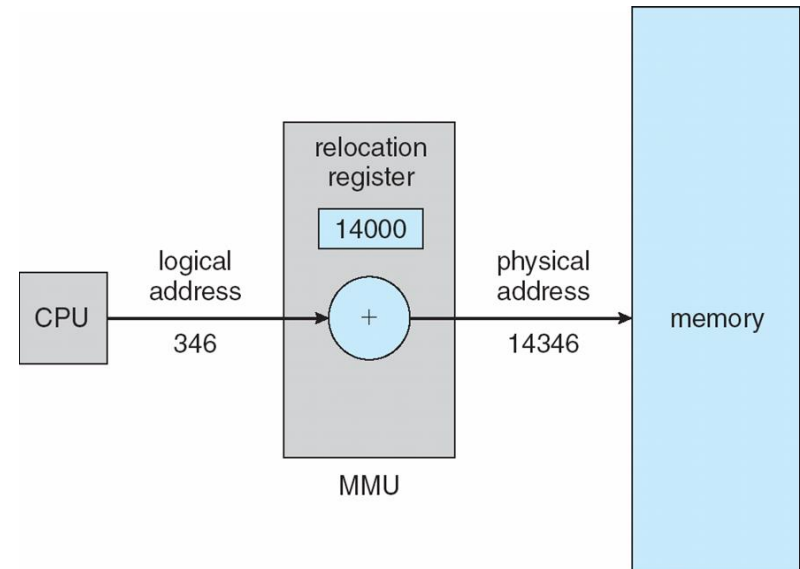- **Physical address space** is the set of all physical addresses generated by a program

# Memory-Management Unit (MMU)

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

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

- To start, consider simple scheme where the value in the relocation register is added to every address generated by a user process at the time it is sent to memory

  - Base register now called **relocation register**

  - MS-DOS on Intel 80x86 used 4 relocation registers

- The user program deals with *logical* addresses; it never sees the *real* physical addresses

  - Execution-time binding occurs when reference is made to location in memory

  - Logical address bound to physical addresses

# Dynamic relocation using a relocation register

- Routine is not loaded until it is called

- Better memory-space utilization; unused routine is never loaded

- All routines kept on disk in relocatable load format

- Useful when large amounts of code are needed to handle infrequently occurring cases

- No special support from the operating system is required

  - Implemented through program design

  - OS can help by providing libraries to implement dynamic loading

# Dynamic Linking

- **Static linking** – system libraries and program code combined by the loader into the binary program image
- Dynamic linking –linking postponed until execution time
- Small piece of code, **stub**, used to locate the appropriate memory-resident library routine
- Stub replaces itself with the address of the routine, and executes the routine
- Operating system checks if routine is in processes' memory address
    - If not in address space, add to address space
- Dynamic linking is particularly useful for libraries
- System also known as **shared libraries**
- Consider applicability to patching system libraries
    - Versioning may be needed

# Swapping

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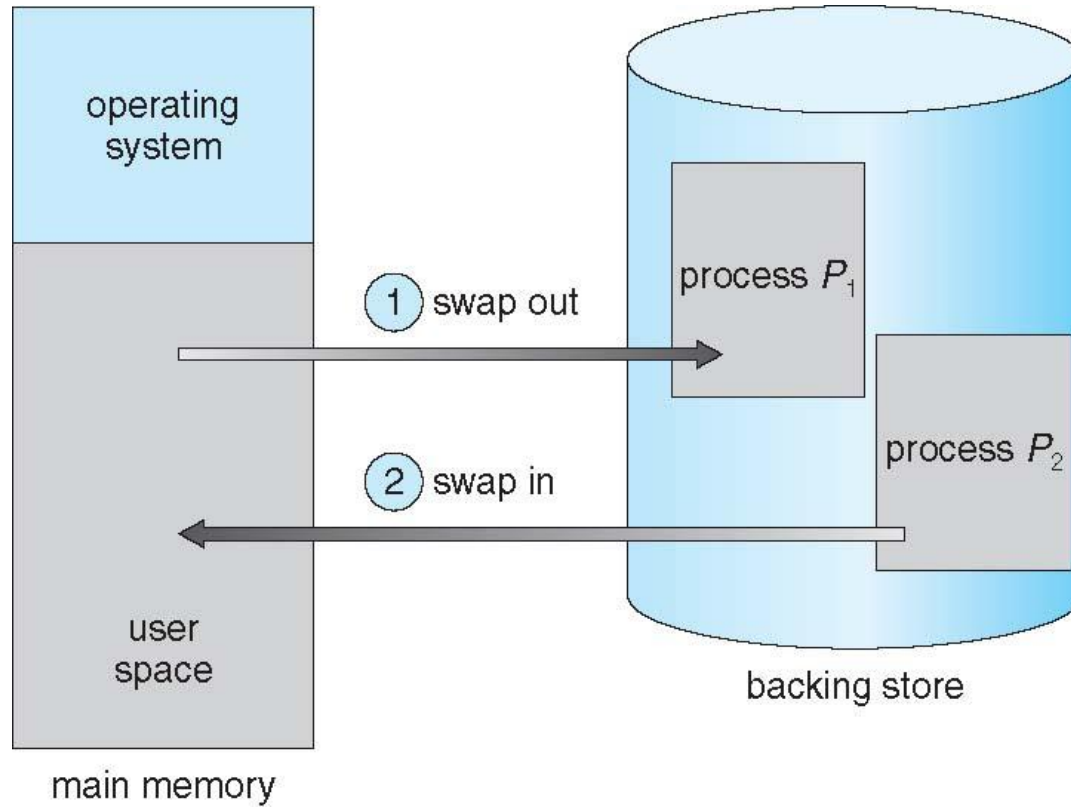- System maintains a **ready queue** of ready-to-run processes which have memory images on disk

# Swapping (Cont.)

- Does the swapped out process need to swap back in to same physical addresses?
- Depends on address binding method
    - Plus consider pending I/O to / from process memory space
- Modified versions of swapping are found on many systems (i.e., UNIX, Linux, and Windows)

    - Swapping normally disabled

    - Started if more than threshold amount of memory allocated

    - Disabled again once memory demand reduced below threshold

# Schematic View of Swapping

# Context Switch Time including Swapping

- If next processes to be put on CPU is not in memory, need to swap out a process and swap in target process

- Context switch time can then be very high

- 100MB process swapping to hard disk with transfer rate of 50MB/sec

  - Swap out time of 2000 ms

  - Plus swap in of same sized process

  - Total context switch swapping component time of 4000ms (4 seconds)

- Can reduce if reduce size of memory swapped – by knowing how much memory really being used

  - System calls to inform OS of memory use via `request_memory()` and `release_memory()`

# Context Switch Time and Swapping (Cont.)

- ☐ Other constraints as well on swapping
  - ☐ Pending I/O – can't swap out as I/O would occur to wrong process
  - ☐ Or always transfer I/O to kernel space, then to I/O device
    - ‣ Known as **double buffering**, adds overhead
- ☐ Standard swapping not used in modern operating systems
  - ☐ But modified version common
    - ‣ Swap only when free memory extremely low

# Swapping on Mobile Systems

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

# Contiguous Allocation

- Main memory must support both OS and user processes

- Limited resource, must allocate efficiently

- Contiguous allocation is one early method

- Main memory usually into two **partitions**:

    - Resident operating system, usually held in low memory with interrupt vector

    - User processes then held in high memory

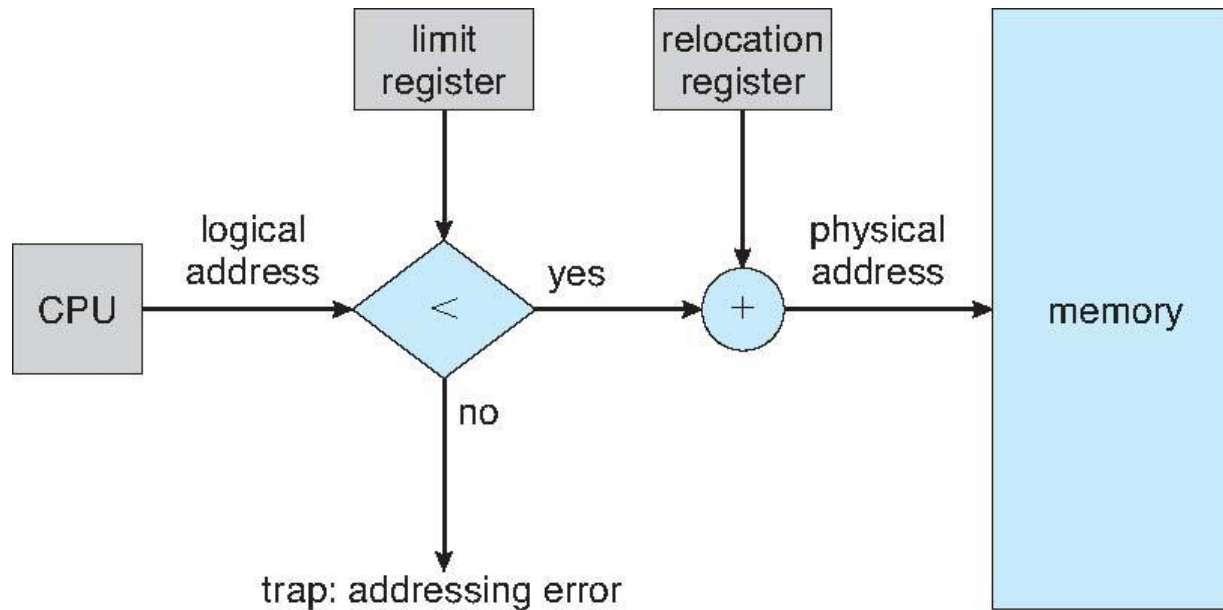    - Each process contained in single contiguous section of memory

# Contiguous Allocation (Cont.)

- Relocation registers used to protect user processes from each other, and from changing operating-system code and data

    - Base register contains value of smallest physical address

    - Limit register contains range of logical addresses – each logical address must be less than the limit register

    - MMU maps logical address *dynamically*

    - Can then allow actions such as kernel code being **transient** and kernel changing size
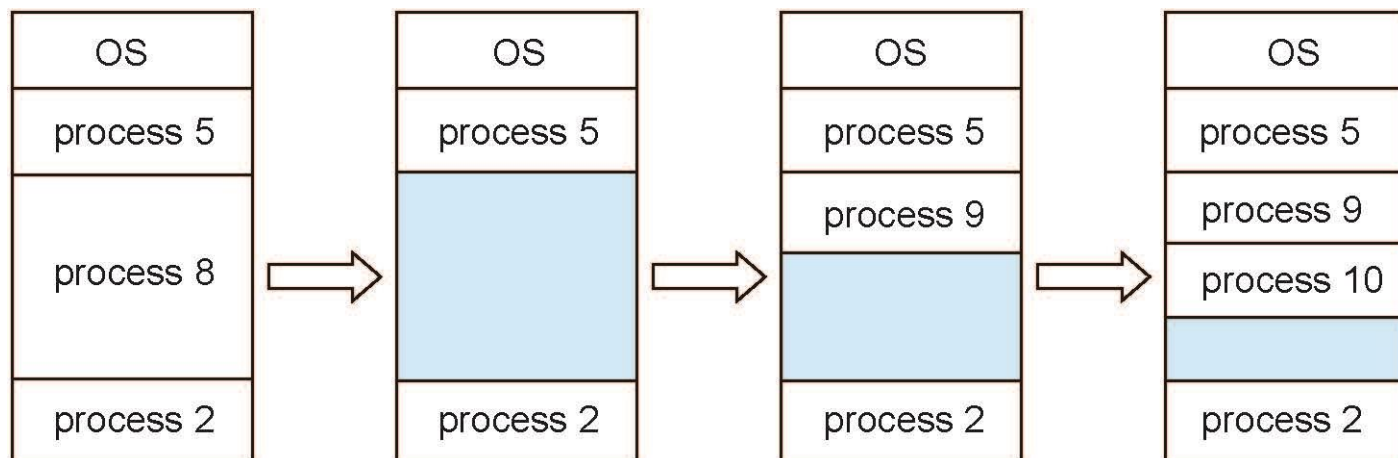
# Multiple-partition allocation

- Multiple-partition allocation
    - Degree of multiprogramming limited by number of partitions
    - **Variable-partition** sizes for efficiency (sized to a given process' needs)
    - **Hole** – block of available memory; holes of various size are scattered throughout memory
    - When a process arrives, it is allocated memory from a hole large enough to accommodate it
    - Process exiting frees its partition, adjacent free partitions combined
    - Operating system maintains information about:
      a) allocated partitions    b) free partitions (hole)

| OS | | OS | | OS | | OS |
|---|---|---|---|---|---|---|
| process 5 | | process 5 | | process 5 | | process 5 |
| | | | | process 9 | | process 9 |
| process 8 | ⇒ | | ⇒ | | ⇒ | process 10 |
| | | | | | | |
| process 2 | | process 2 | | process 2 | | process 2 |

# Dynamic Storage-Allocation Problem

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

- **First-fit**:  Allocate the *first* hole that is big enough

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

- **Worst-fit**:  Allocate the *largest* hole; must also search entire list
  - Produces the largest leftover hole

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

# Fragmentation

- **External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous

- **Internal Fragmentation** – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used

- First fit analysis reveals that given $N$ blocks allocated, $0.5 N$ blocks lost to fragmentation

  - 1/3 may be unusable -> **50-percent rule**

# Fragmentation (Cont.)

- Reduce external fragmentation by **compaction**

    - Shuffle memory contents to place all free memory together in one large block

    - Compaction is possible *only* if relocation is dynamic, and is done at execution time

    - I/O problem

        ▸ Latch job in memory while it is involved in I/O

        ▸ Do I/O only into OS buffers

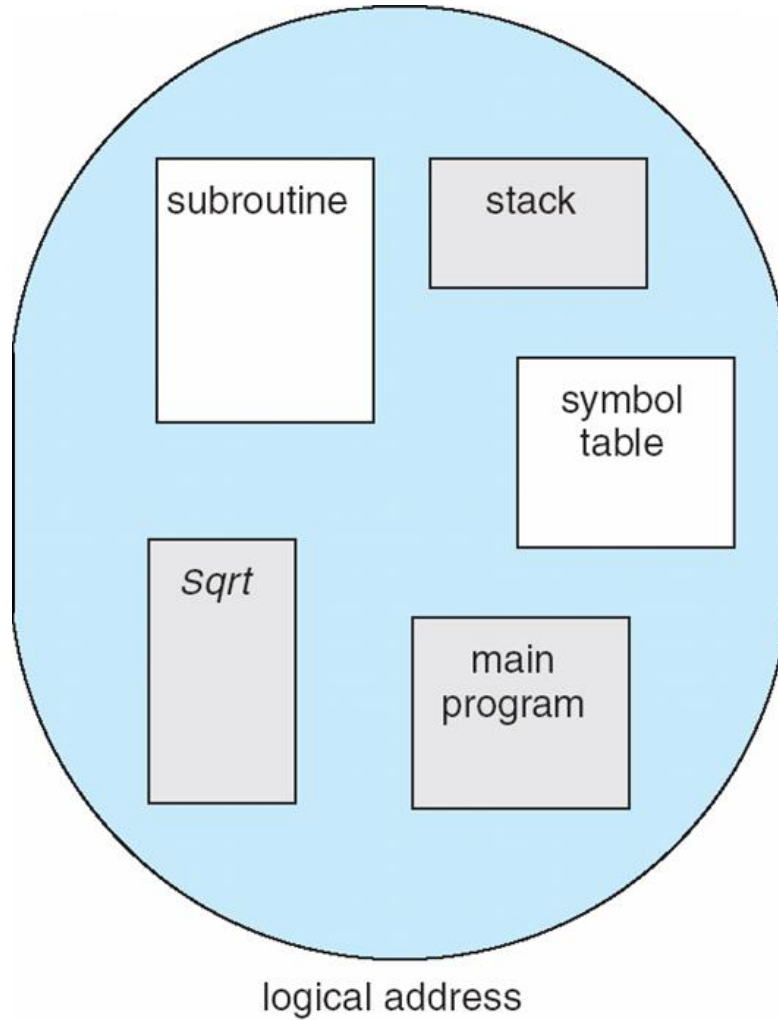- Now consider that backing store has same fragmentation problems

# Segmentation

- Memory-management scheme that supports user view of memory
- A program is a collection of segments
    - A segment is a logical unit such as:

        main program

        procedure

        function

        method

        object

        local variables, global variables

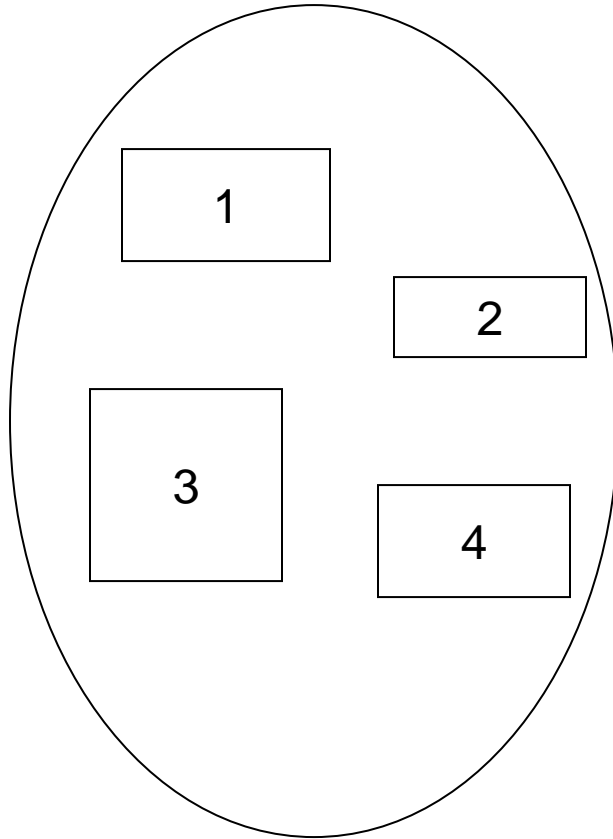        common block

        stack

        symbol table
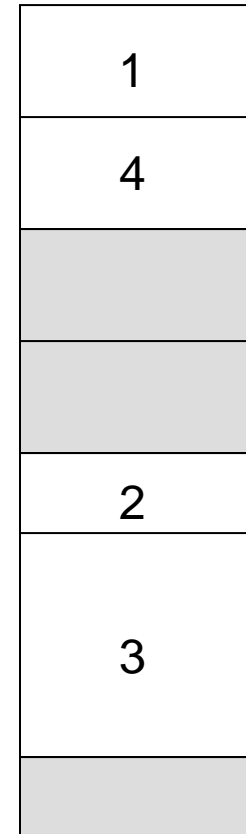
        arrays

# User's View of a Program

# Logical View of Segmentation

user space

physical memory space

# Segmentation Architecture

- Logical address consists of a two tuple:

   <segment-number, offset>,

- **Segment table** – maps two-dimensional physical addresses; each table entry has:

   - **base** – contains the starting physical address where the segments reside in memory
   - **limit** – specifies the length of the segment

- **Segment-table base register (STBR)** points to the segment table's location in memory

- **Segment-table length register (STLR)** indicates number of segments used by a program;

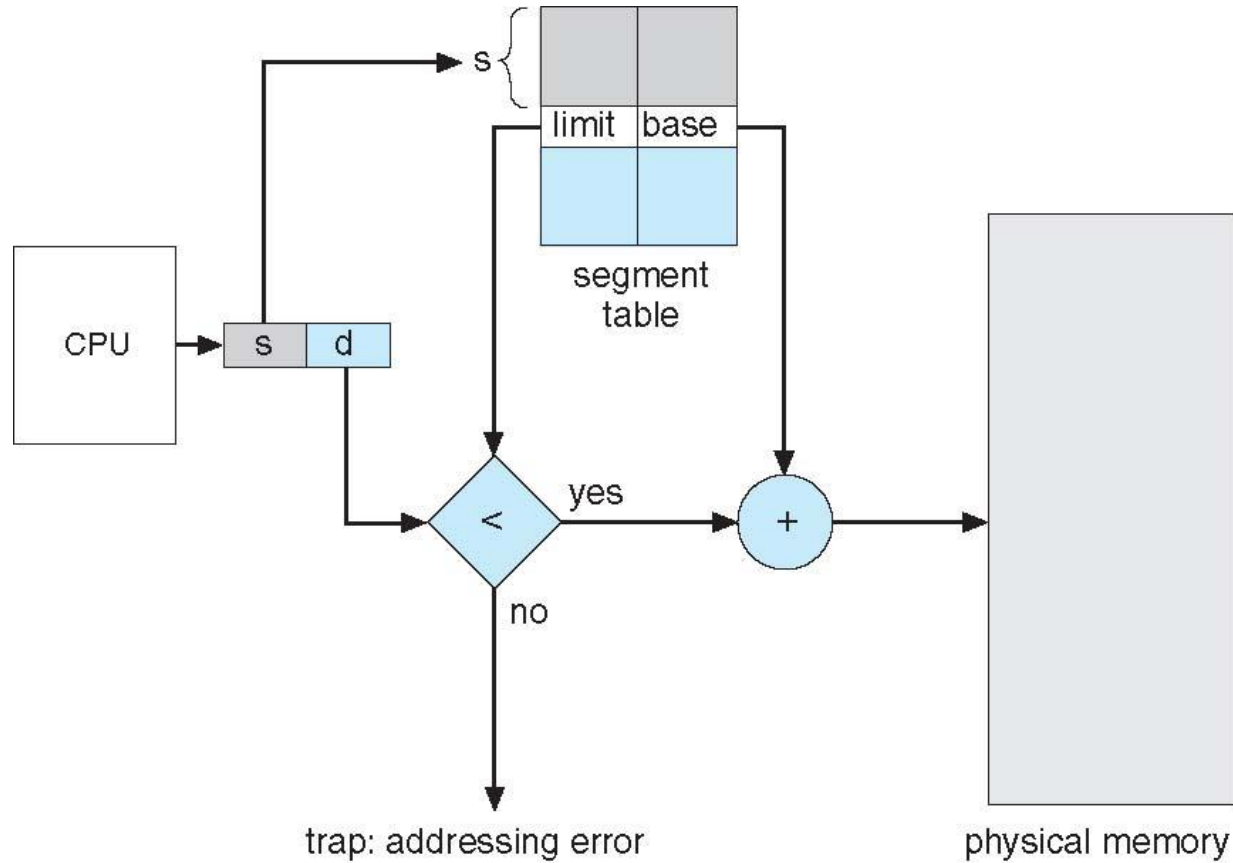   segment number **s** is legal if **s** < **STLR**

# Segmentation Architecture (Cont.)

- Protection
    - With each entry in segment table associate:
        - validation bit = 0 $\Rightarrow$ illegal segment
        - read/write/execute privileges
- Protection bits associated with segments; code sharing occurs at segment level
- Since segments vary in length, memory allocation is a dynamic storage-allocation problem
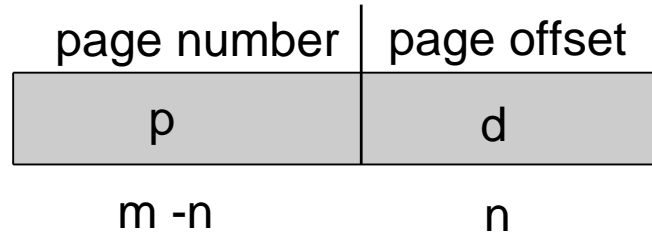- A segmentation example is shown in the following diagram

# Segmentation Hardware

# Paging

- Physical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available
    - Avoids external fragmentation
    - Avoids problem of varying sized memory chunks
- Divide physical memory into fixed-sized blocks called **frames**
    - Size is power of 2, between 512 bytes and 16 Mbytes
- Divide logical memory into blocks of same size called **pages**
- Keep track of all free frames
- To run a program of size *N* pages, need to find *N* free frames and load program
- Set up a **page table** to translate logical to physical addresses
- Backing store likewise split into pages
- Still have Internal fragmentation
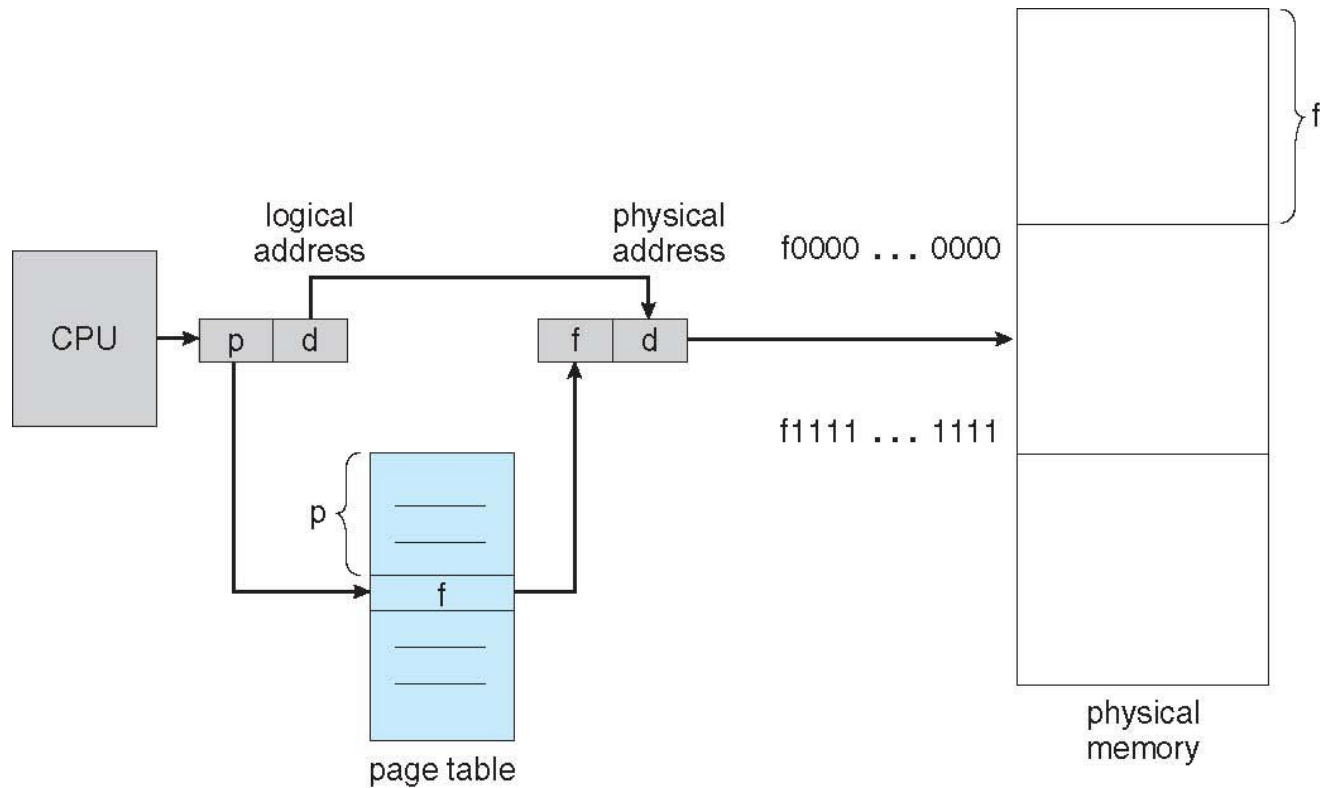
# Address Translation Scheme

- Address generated by CPU is divided into:

  - **Page number** (*p*) – used as an index into a **page table** which contains base address of each page in physical memory

  - **Page offset** (*d*) – combined with base address to define the physical memory address that is sent to the memory unit

| page number | page offset |
|:---:|:---:|
| p | d |
| m -n | n |

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

# Paging Hardware

# Paging Model of Logical and Physical Memory

# Paging Example



n=2 and m=4   32-byte memory and 4-byte pages

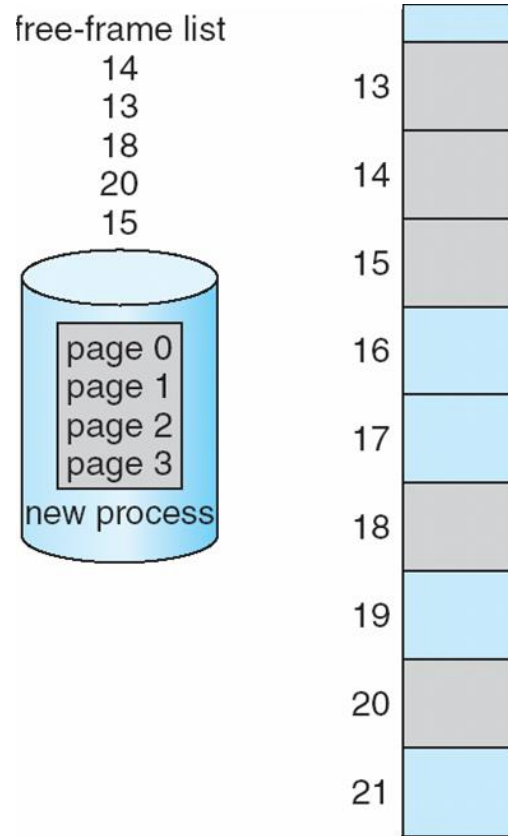# Paging (Cont.)

- Calculating internal fragmentation
    - Page size = 2,048 bytes
    - Process size = 72,766 bytes
    - 35 pages + 1,086 bytes
    - Internal fragmentation of 2,048 - 1,086 = 962 bytes
    - Worst case fragmentation = 1 frame – 1 byte
    - On average fragmentation = 1 / 2 frame size
    - So small frame sizes desirable?
    - But each page table entry takes memory to track
    - Page sizes growing over time
        - Solaris supports two page sizes – 8 KB and 4 MB
- Process view and physical memory now very different
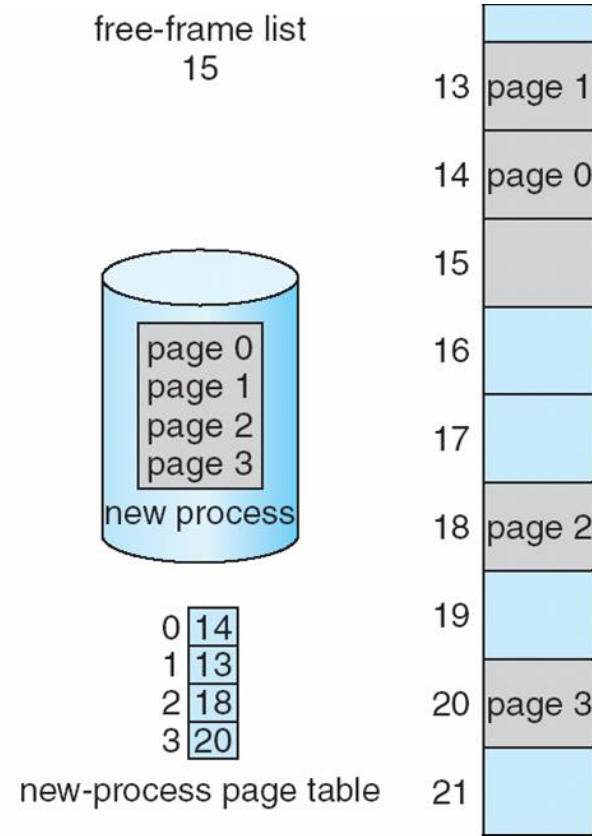- By implementation process can only access its own memory

# Free Frames



Before allocation                           After allocation

# Implementation of Page Table

- Page table is kept in main memory

- **Page-table base register** (**PTBR**) points to the page table

- **Page-table length register** (**PTLR**) indicates size of the page table

- In this scheme every data/instruction access requires two memory accesses

    - One for the page table and one for the data / instruction

- The two memory access problem can be solved by the use of a special fast-lookup hardware cache called **associative memory** or **translation look-aside buffers** (**TLBs**)

# Implementation of Page Table (Cont.)

- Some TLBs store **address-space identifiers** (**ASIDs**) in each TLB entry – uniquely identifies each process to provide address-space protection for that process

    - Otherwise need to flush at every context switch

- TLBs typically small (64 to 1,024 entries)

- On a TLB miss, value is loaded into the TLB for faster access next time

    - Replacement policies must be considered

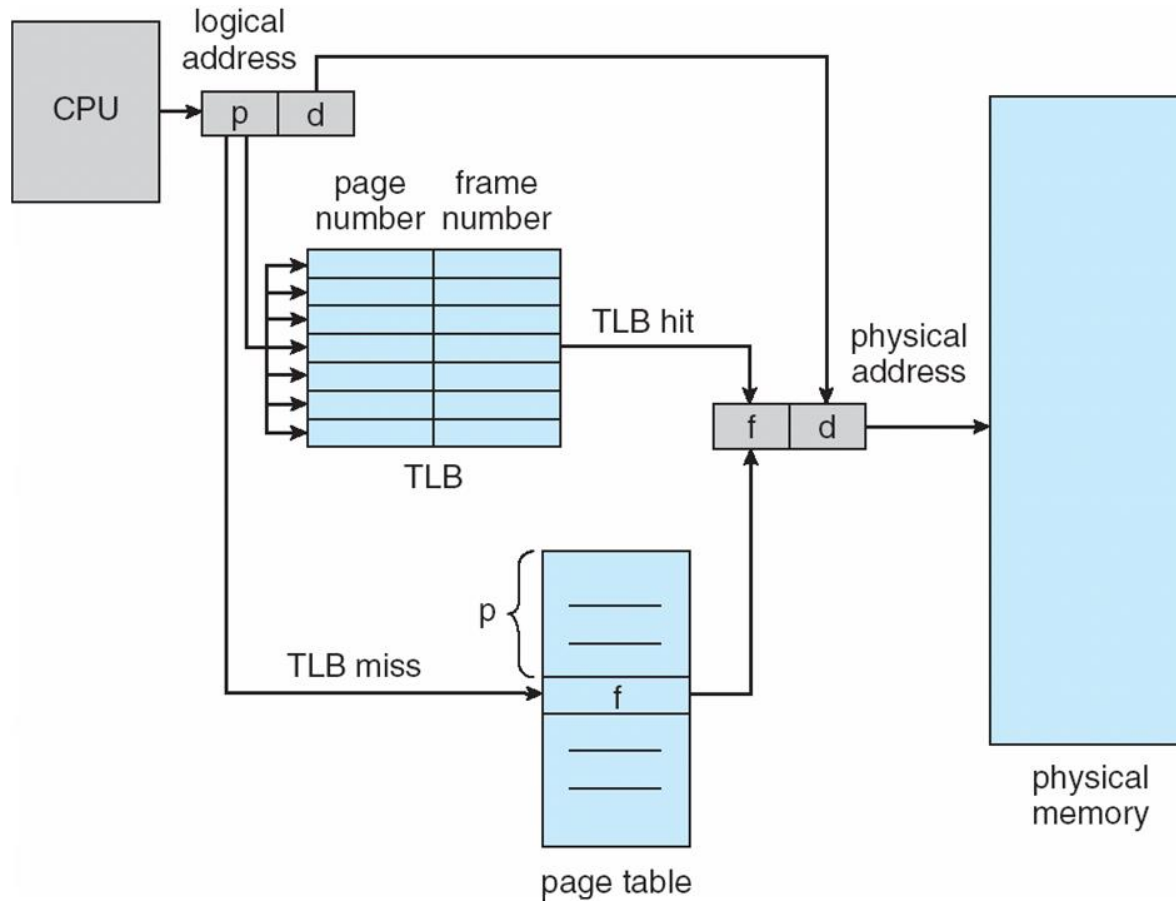    - Some entries can be **wired down** for permanent fast access

# Associative Memory

□ Associative memory – parallel search

| Page # | Frame # |
|--------|---------|
|        |         |
|        |         |
|        |         |
|        |         |

□ Address translation (p, d)

 □ If p is in associative register, get frame # out

 □ Otherwise get frame # from page table in memory

# Paging Hardware With TLB

# Effective Access Time

- Associative Lookup = $\varepsilon$ time unit
  - Can be < 10% of memory access time
- Hit ratio = $\alpha$
  - Hit ratio – percentage of times that a page number is found in the associative registers; ratio related to number of associative registers
- Consider $\alpha$ = 80%, $\varepsilon$ = 20ns for TLB search, 100ns for memory access
- **Effective Access Time** (**EAT**)

$$EAT = (1 + \varepsilon)\,\alpha + (2 + \varepsilon)(1 - \alpha)$$
$$= 2 + \varepsilon - \alpha$$

- Consider $\alpha$ = 80%, $\varepsilon$ = 20ns for TLB search, 100ns for memory access
  - EAT = 0.80 x 100 + 0.20 x 200 = 120ns
- Consider more realistic hit ratio -> $\alpha$ = 99%, $\varepsilon$ = 20ns for TLB search, 100ns for memory access
  - EAT = 0.99 x 100 + 0.01 x 200 = 101ns

# Memory Protection

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

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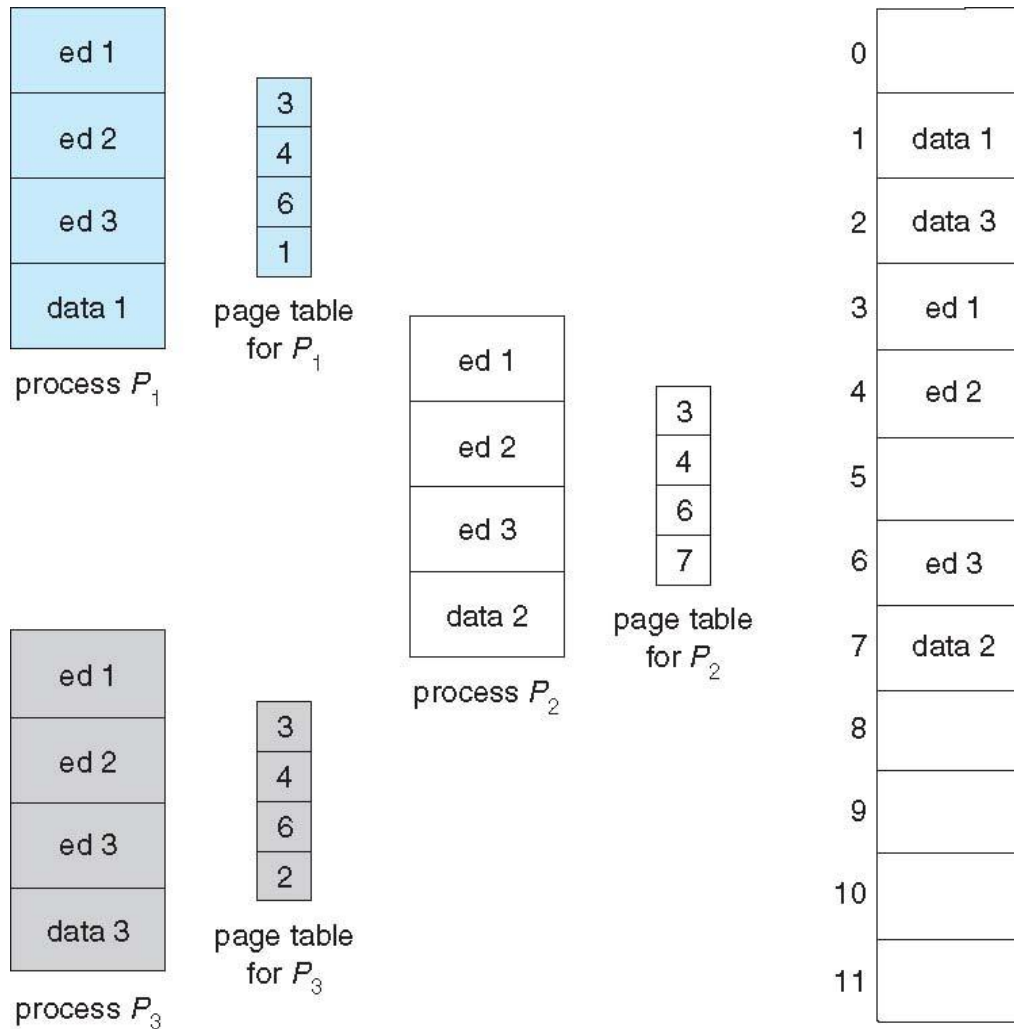
# Shared Pages

- **Shared code**
    - One copy of read-only (**reentrant**) code shared among processes (i.e., text editors, compilers, window systems)
    - Similar to multiple threads sharing the same process space
    - Also useful for interprocess communication if sharing of read-write pages is allowed

- **Private code and data**
    - Each process keeps a separate copy of the code and data
    - The pages for the private code and data can appear anywhere in the logical address space

# Shared Pages Example

# Structure of the Page Table

- Memory structures for paging can get huge using straight-forward methods

  - Consider a 32-bit logical address space as on modern computers

  - Page size of 4 KB ($2^{12}$)

  - Page table would have 1 million entries ($2^{32}$ / $2^{12}$)

  - If each entry is 4 bytes -> 4 MB of physical address space / memory for page table alone

    - That amount of memory used to cost a lot

    - Don't want to allocate that contiguously in main memory

- Hierarchical Paging

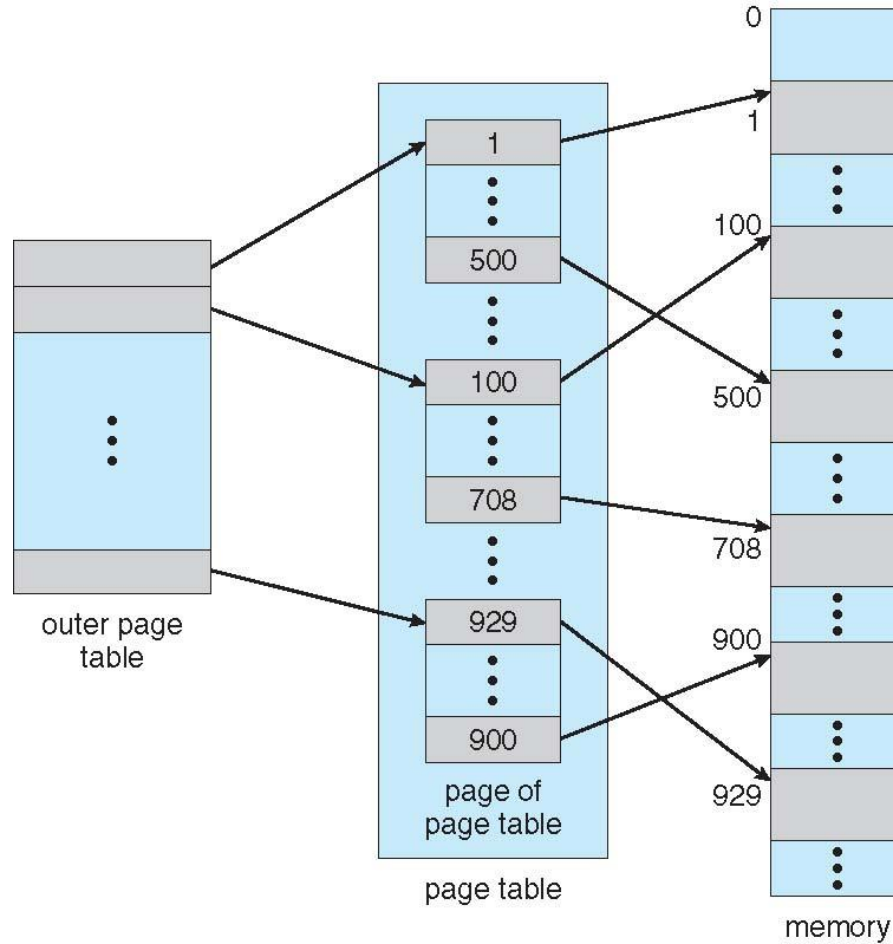- Hashed Page Tables

- Inverted Page Tables

# Hierarchical Page Tables

- Break up the logical address space into multiple page tables

- A simple technique is a two-level page table

- We then page the page table

# Two-Level Page-Table Scheme
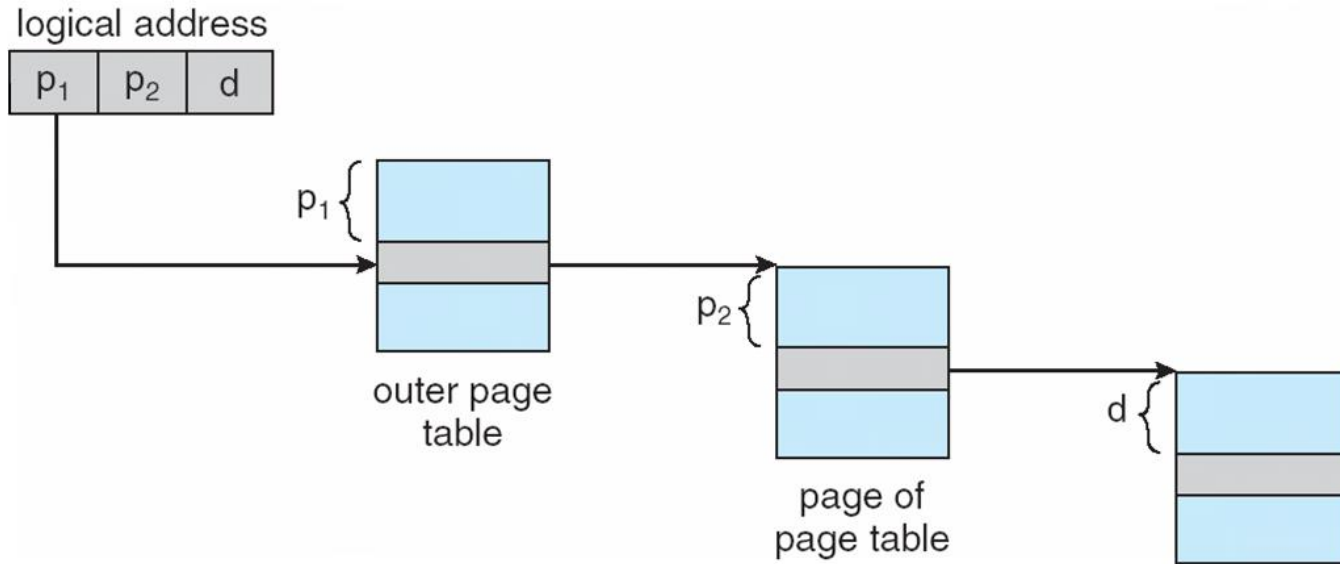
# Two-Level Paging Example

- A logical address (on 32-bit machine with 1K page size) is divided into:
  - a page number consisting of 22 bits
  - a page offset consisting of 10 bits

- Since the page table is paged, the page number is further divided into:
  - a 12-bit page number
  - a 10-bit page offset

- Thus, a logical address is as follows:

| page number | | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 12 | 10 | 10 |

- where $p_1$ is an index into the outer page table, and $p_2$ is the displacement within the page of the inner page table
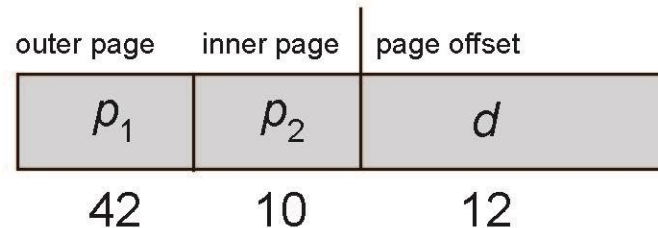- Known as **forward-mapped page table**

# Address-Translation Scheme



logical address

| $p_1$ | $p_2$ | d |

$p_1$ { outer page table

$p_2$ { page of page table

d {

# 64-bit Logical Address Space

- Even two-level paging scheme not sufficient
- If page size is 4 KB ($2^{12}$)
  - Then page table has $2^{52}$ entries
  - If two level scheme, inner page tables could be $2^{10}$ 4-byte entries
  - Address would look like

| outer page | inner page | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

  - Outer page table has $2^{42}$ entries or $2^{44}$ bytes
  - One solution is to add a 2nd outer page table
  - But in the following example the 2nd outer page table is still $2^{34}$ bytes in size
    - ▸ And possibly 4 memory access to get to one physical memory location

# Three-level Paging Scheme

| outer page | inner page | offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

| 2nd outer page | outer page | inner page | offset |
|:---:|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $p_3$ | $d$ |
| 32 | 10 | 10 | 12 |

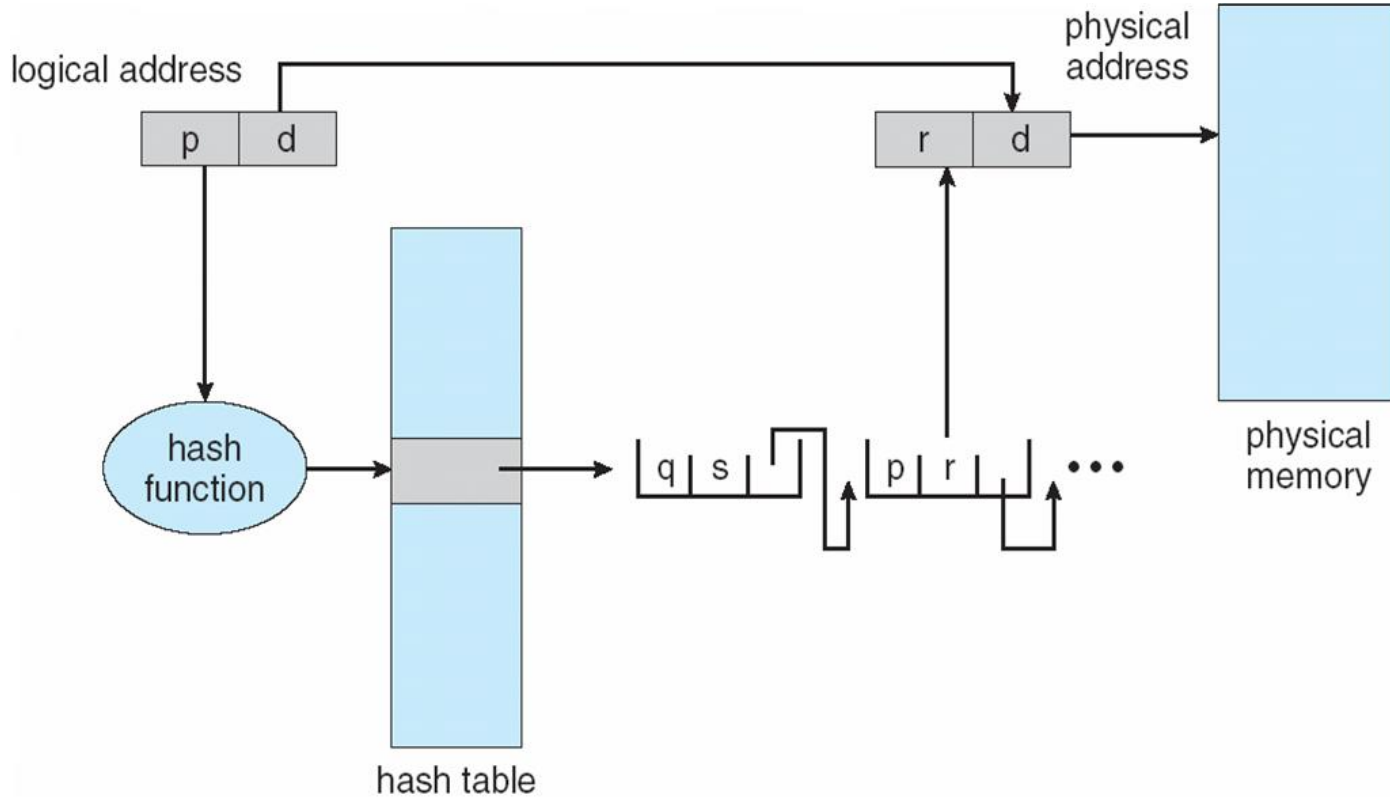# Hashed Page Tables

- Common in address spaces > 32 bits

- The virtual page number is hashed into a page table

    - This page table contains a chain of elements hashing to the same location

- Each element contains (1) the virtual page number (2) the value of the mapped page frame (3) a pointer to the next element

- Virtual page numbers are compared in this chain searching for a match

    - If a match is found, the corresponding physical frame is extracted

- Variation for 64-bit addresses is **clustered page tables**

    - Similar to hashed but each entry refers to several pages (such as 16) rather than 1

    - Especially useful for **sparse** address spaces (where memory references are non-contiguous and scattered)
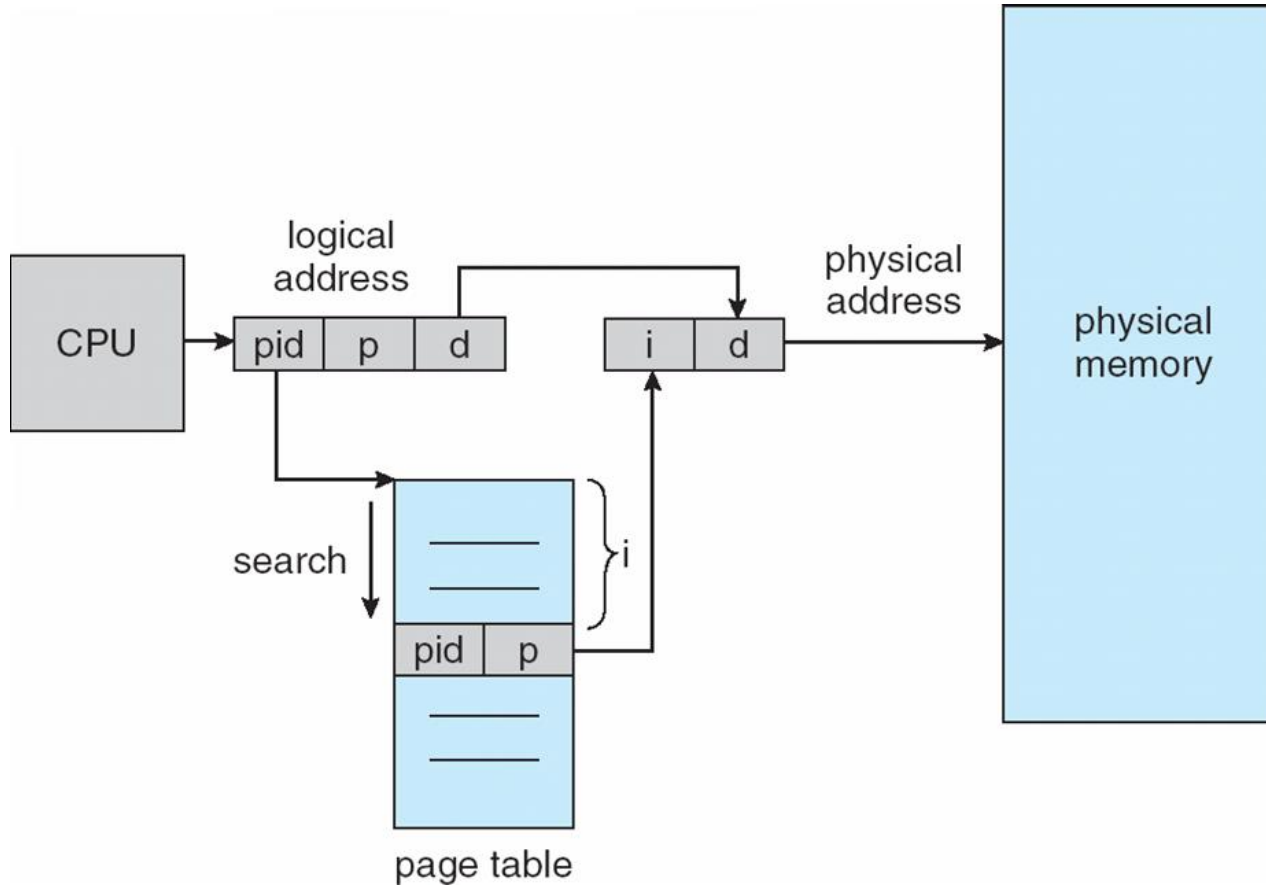
# Hashed Page Table

# Inverted Page Table

- Rather than each process having a page table and keeping track of all possible logical pages, track all physical pages

- One entry for each real page of memory

- Entry consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page

- Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs

- Use hash table to limit the search to one — or at most a few — page-table entries
  - TLB can accelerate access

- But how to implement shared memory?
  - One mapping of a virtual address to the shared physical address

# Inverted Page Table Architecture

# Oracle SPARC Solaris

- Consider modern, 64-bit operating system example with tightly integrated HW
  - Goals are efficiency, low overhead
- Based on hashing, but more complex
- Two hash tables
  - One kernel and one for all user processes
  - Each maps memory addresses from virtual to physical memory
  - Each entry represents a contiguous area of mapped virtual memory,
    - More efficient than having a separate hash-table entry for each page
  - Each entry has base address and span (indicating the number of pages the entry represents)

# Oracle SPARC Solaris (Cont.)

- TLB holds translation table entries (TTEs) for fast hardware lookups
  - A cache of TTEs reside in a translation storage buffer (TSB)
    - Includes an entry per recently accessed page
- Virtual address reference causes TLB search
  - If miss, hardware walks the in-memory TSB looking for the TTE corresponding to the address
    - If match found, the CPU copies the TSB entry into the TLB and translation completes
    - If no match found, kernel interrupted to search the hash table
      - The kernel then creates a TTE from the appropriate hash table and stores it in the TSB, Interrupt handler returns control to the MMU, which completes the address translation.

# Example: The Intel 32 and 64-bit Architectures

☐ Dominant industry chips

☐ Pentium CPUs are 32-bit and called IA-32 architecture

☐ Current Intel CPUs are 64-bit and called IA-64 architecture

☐ Many variations in the chips, cover the main ideas here
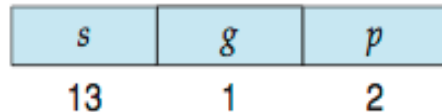
# Example: The Intel IA-32 Architecture

- Supports both segmentation and segmentation with paging

  - Each segment can be 4 GB

  - Up to 16 K segments per process

  - Divided into two partitions

    - First partition of up to 8 K segments are private to process (kept in **local descriptor table** (**LDT**))

    - Second partition of up to 8K segments shared among all processes (kept in **global descriptor table** (**GDT**))

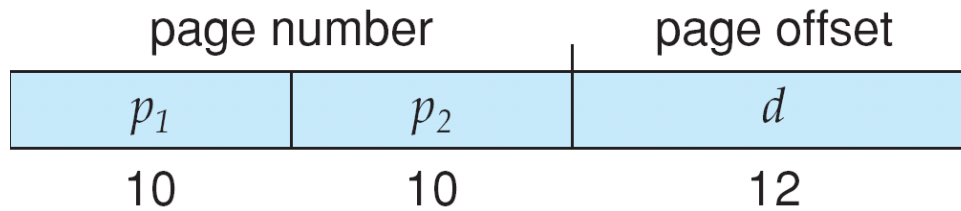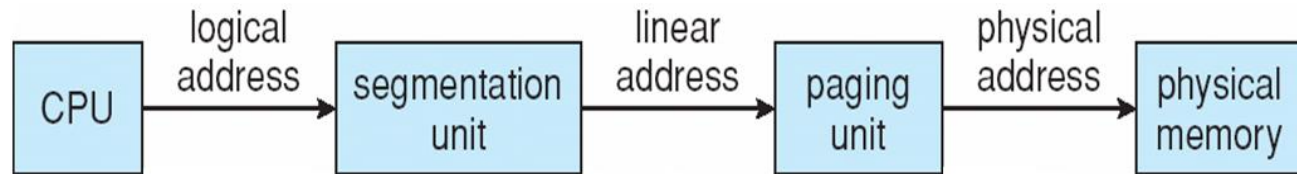# Example: The Intel IA-32 Architecture (Cont.)

- CPU generates logical address
    - Selector given to segmentation unit
        - ▸ Which produces linear addresses

| $s$ | $g$ | $p$ |
|:---:|:---:|:---:|
| 13 | 1 | 2 |

    - Linear address given to paging unit
        - ▸ Which generates physical address in main memory
        - ▸ Paging units form equivalent of MMU
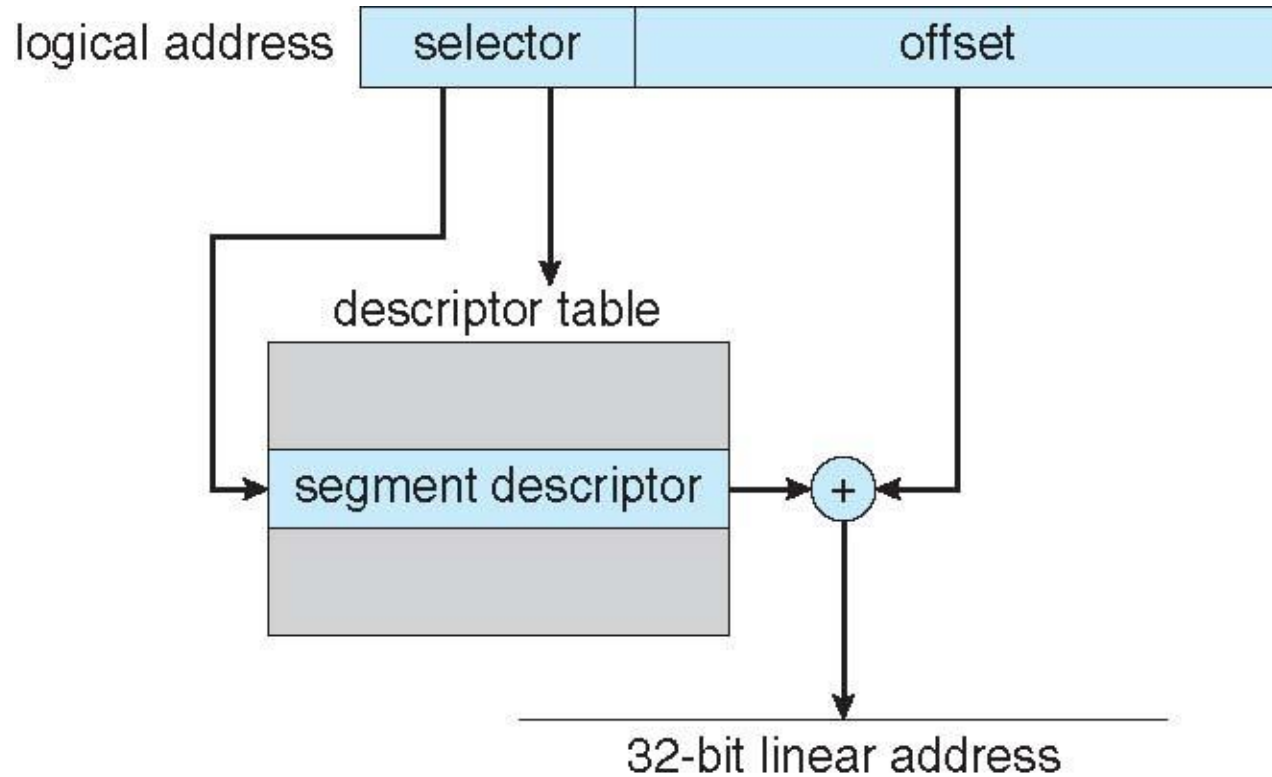        - ▸ Pages sizes can be 4 KB or 4 MB

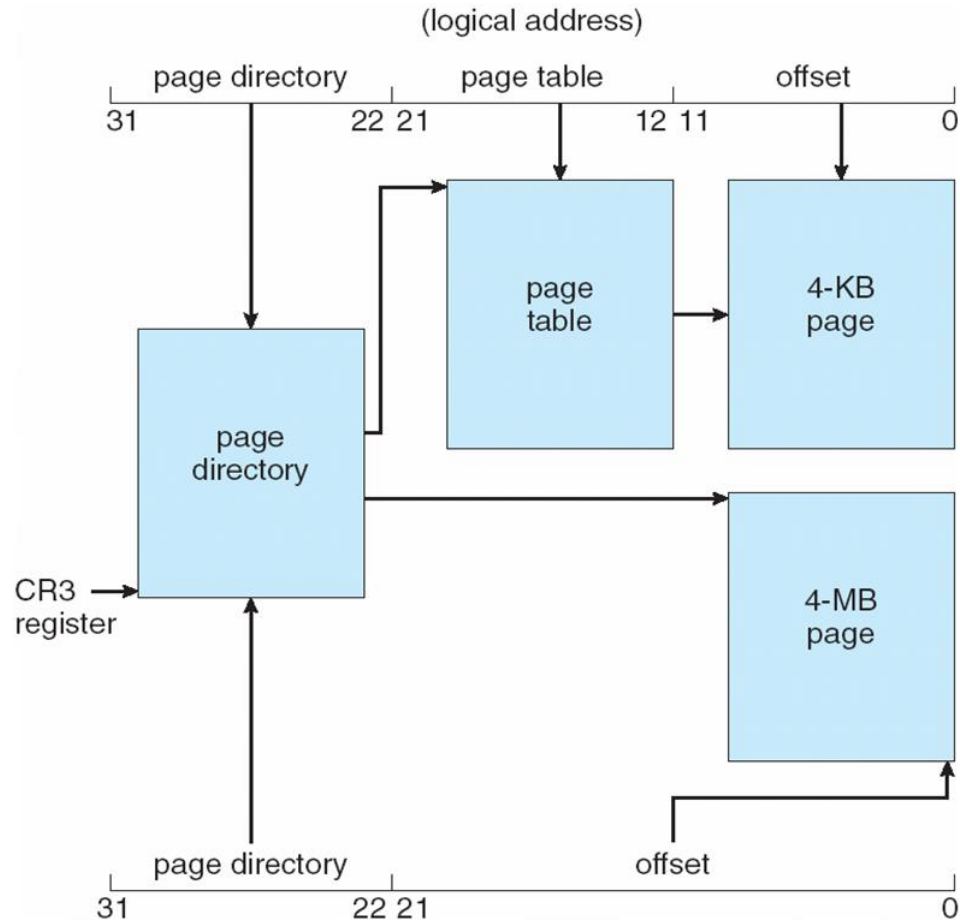# Logical to Physical Address Translation in IA-32
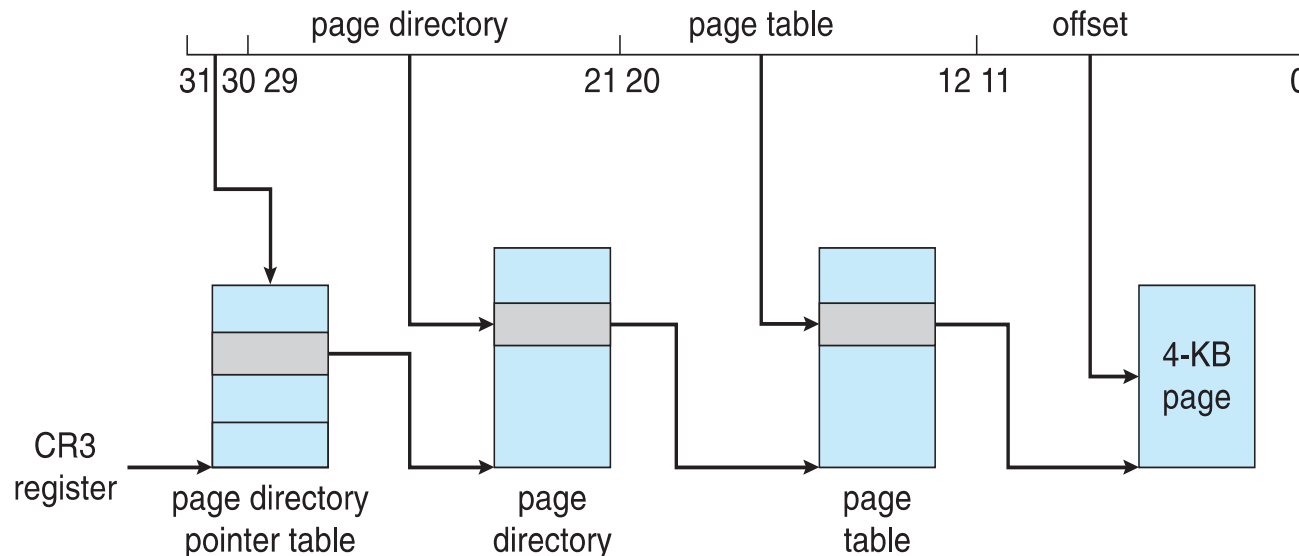
# Intel IA-32 Segmentation

# Intel IA-32 Paging Architecture

# Intel IA-32 Page Address Extensions

- 32-bit address limits led Intel to create **page address extension** (**PAE**), allowing 32-bit apps access to more than 4GB of memory space
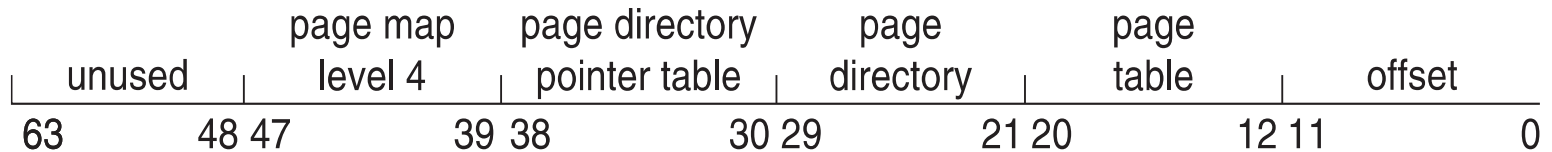
  - Paging went to a 3-level scheme

  - Top two bits refer to a **page directory pointer table**

  - Page-directory and page-table entries moved to 64-bits in size

  - Net effect is increasing address space to 36 bits – 64GB of physical memory

# Intel x86-64

- Current generation Intel x86 architecture

- 64 bits is ginormous (> 16 exabytes)

- In practice only implement 48 bit addressing

    - Page sizes of 4 KB, 2 MB, 1 GB

    - Four levels of paging hierarchy

- Can also use PAE so virtual addresses are 48 bits and physical addresses are 52 bits

| unused | page map level 4 | page directory pointer table | page directory | page table | offset |
|---|---|---|---|---|---|
| 63          48 | 47          39 | 38          30 | 29          21 | 20          12 | 11          0 |

# Example: ARM Architecture

- Dominant mobile platform chip (Apple iOS and Google Android devices for example)

- Modern, energy efficient, 32-bit CPU

- 4 KB and 16 KB pages

- 1 MB and 16 MB pages (termed **sections**)

- One-level paging for sections, two-level for smaller pages

- Two levels of TLBs
  - Outer level has two micro TLBs (one data, one instruction)
  - Inner is single main TLB
  - First inner is checked, on miss outers are checked, and on miss page table walk performed by CPU

| 32 bits | | |
|---|---|---|
| outer page | inner page | offset |

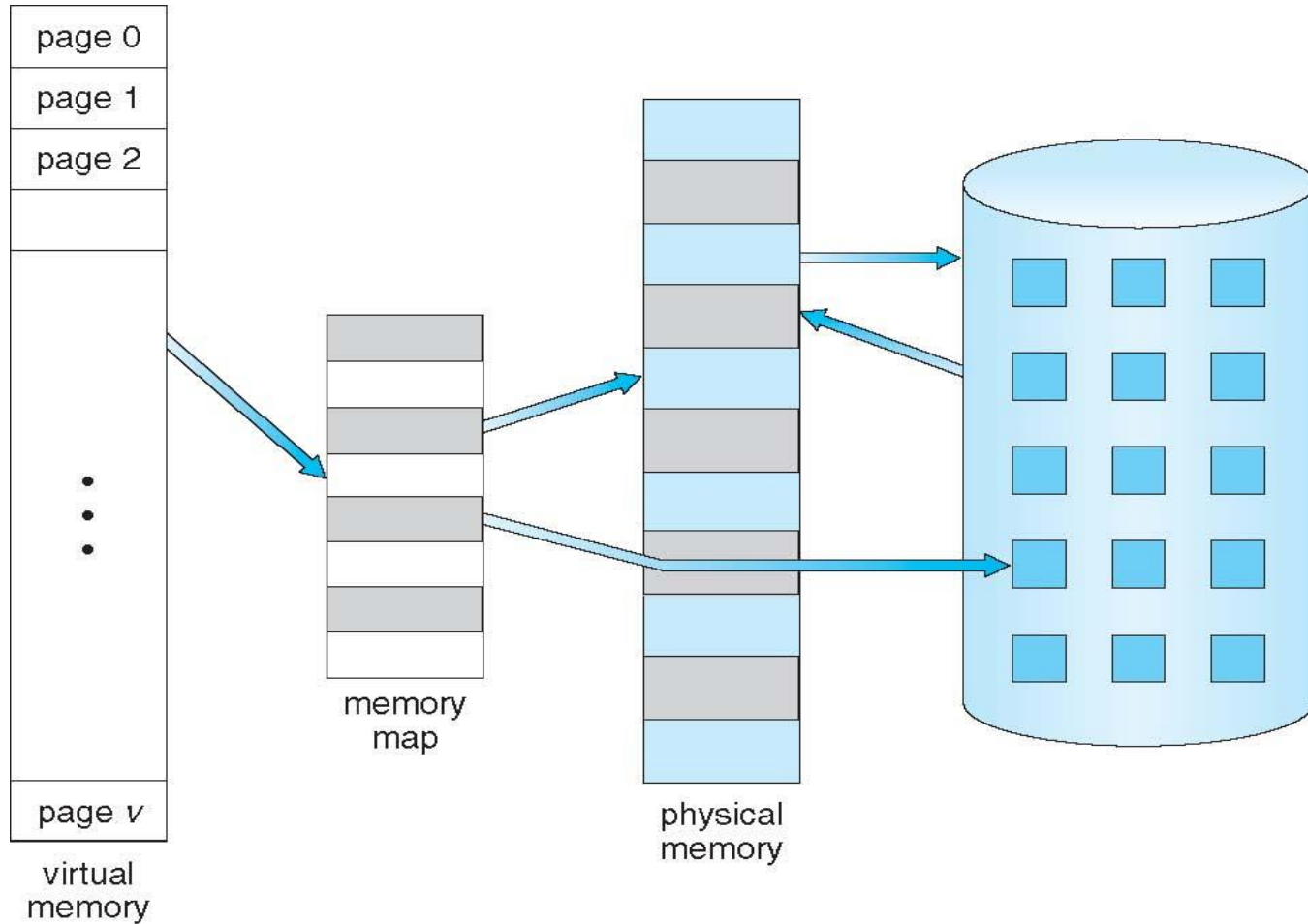- 4-KB or 16-KB page
- 1-MB or 16-MB section

# Virtual Memory

☐ **Virtual memory** – separation of user logical memory from physical memory

☐ Only part of the program needs to be in memory for execution

☐ Logical address space can therefore be much larger than physical address space

☐ Allows address spaces to be shared by several processes

☐ Allows for more efficient process creation

☐ More programs running concurrently

☐ Less I/O needed to load or swap processes

☐ Virtual memory can be implemented via:
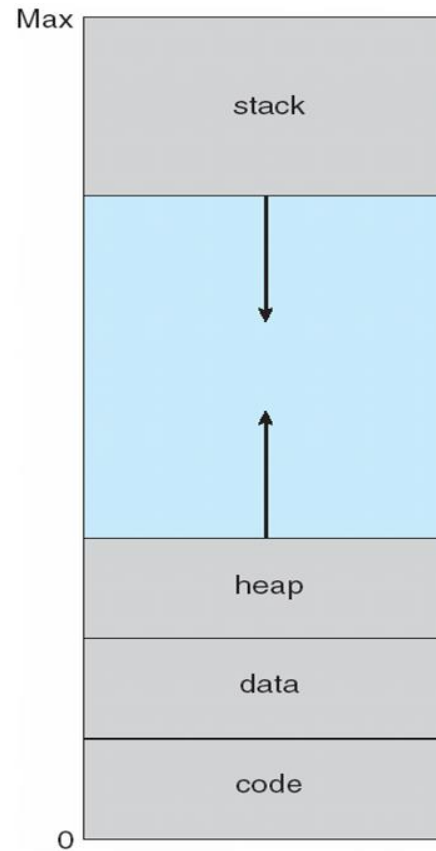
☐ Demand paging

☐ Demand segmentation

# Virtual Memory That is Larger Than Physical Memory



page 0
page 1
page 2

page *v*

virtual memory

memory map

physical memory

# Virtual-address Space

# Virtual Address Space

- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc

- System libraries shared via mapping into virtual address space

- Shared memory by mapping pages read-write into virtual address space

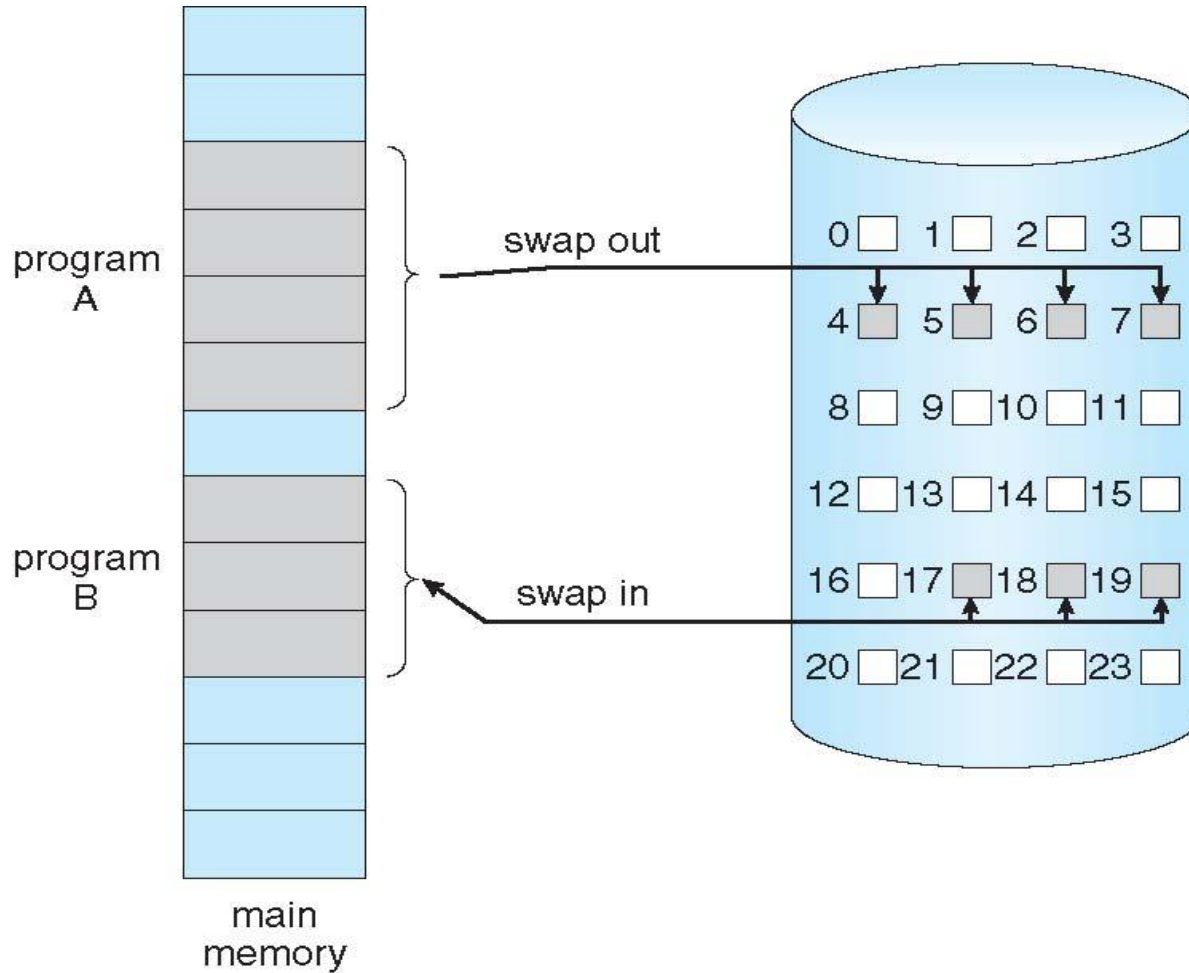- Pages can be shared during `fork()`, speeding process creation

# Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed
  - Less I/O needed, no unnecessary I/O
  - Less memory needed
  - Faster response
  - More users

- Page is needed $\Rightarrow$ reference to it
  - invalid reference $\Rightarrow$ abort
  - not-in-memory $\Rightarrow$ bring to memory

- **Lazy swapper** – never swaps a page into memory unless page will be needed
  - Swapper that deals with pages is a **pager**

# Transfer of a Paged Memory to Contiguous Disk Space

# Valid-Invalid Bit

☐ With each page table entry a valid–invalid bit is associated
(**v** $\Rightarrow$ in-memory – **memory resident**, **i** $\Rightarrow$ not-in-memory)

☐ Initially valid–invalid bit is set to **i** on all entries

☐ Example of a page table snapshot:

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

page table

☐ During address translation, if valid–invalid bit in page table entry
is **I** $\Rightarrow$ page fault

# Page Table When Some Pages Are Not in Main Memory

# Page Fault

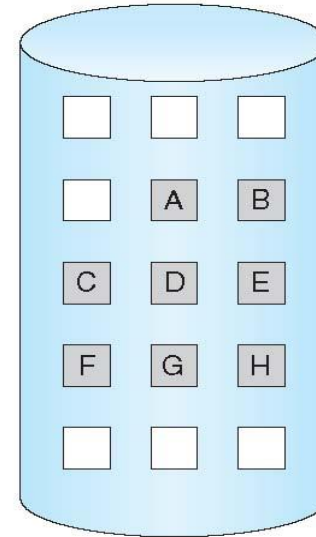- If there is a reference to a page, first reference to that page will trap to operating system:

    **page fault**

1. Operating system looks at another table to decide:
    - Invalid reference $\Rightarrow$ abort
    - Just not in memory

2. Get empty frame

3. Swap page into frame via scheduled disk operation

4. Reset tables to indicate page now in memory
   Set validation bit = **v**

5. Restart the instruction that caused the page fault
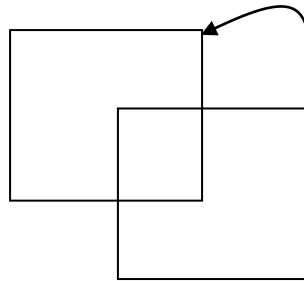
# Aspects of Demand Paging

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

# Instruction Restart

- Consider an instruction that could access several different locations
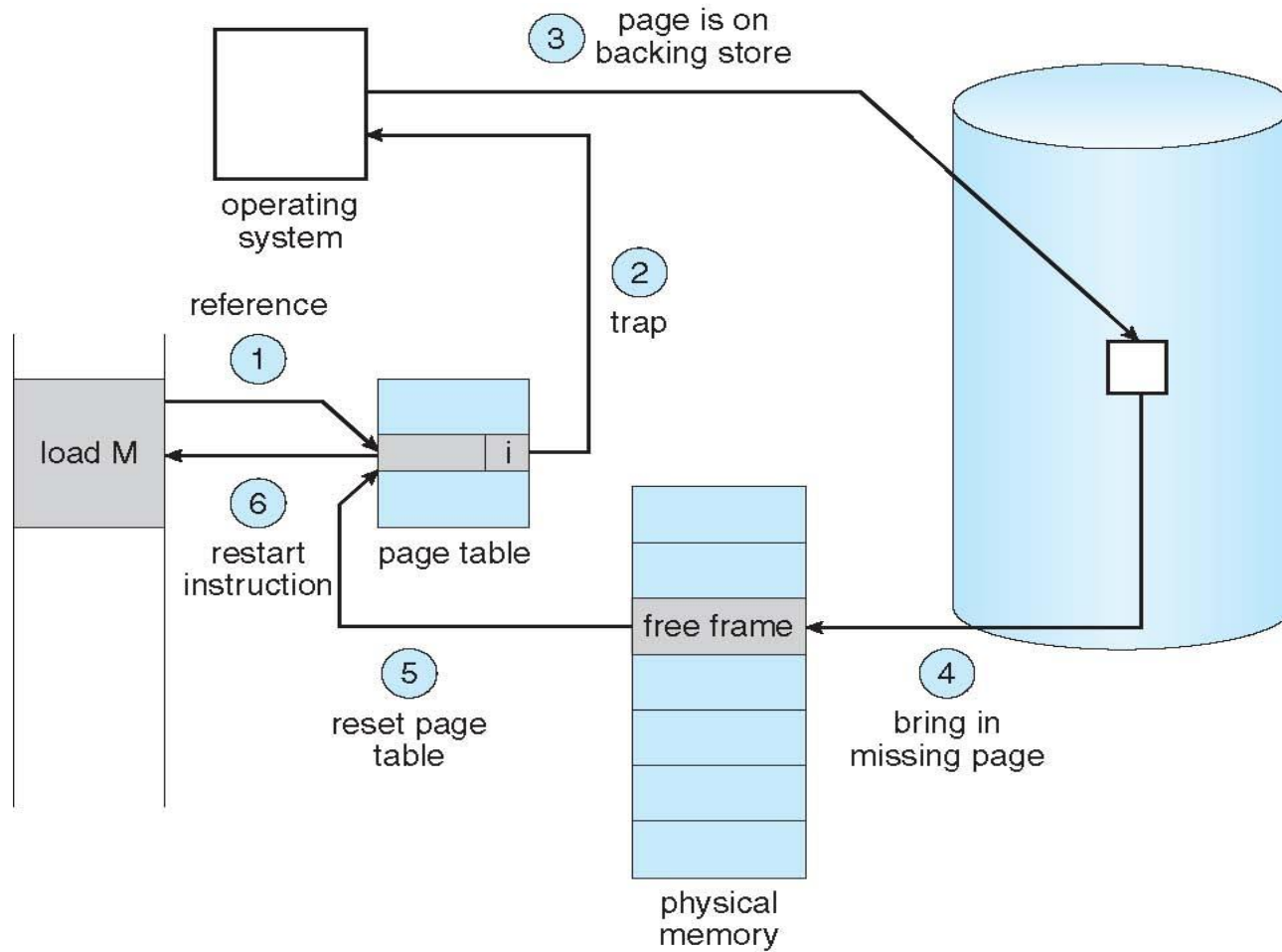  - block move

  

  - auto increment/decrement location
  - Restart the whole operation?
    - What if source and destination overlap?

# Steps in Handling a Page Fault

# Performance of Demand Paging

- Stages in Demand Paging
1. Trap to the operating system
2. Save the user registers and process state
3. Determine that the interrupt was a page fault
4. Check that the page reference was legal and determine the location of the page on the disk
5. Issue a read from the disk to a free frame:
    1. Wait in a queue for this device until the read request is serviced
    2. Wait for the device seek and/or latency time
    3. Begin the transfer of the page to a free frame
6. While waiting, allocate the CPU to some other user
7. Receive an interrupt from the disk I/O subsystem (I/O completed)
8. Save the registers and process state for the other user
9. Determine that the interrupt was from the disk
10. Correct the page table and other tables to show page is now in memory
11. Wait for the CPU to be allocated to this process again
12. Restore the user registers, process state, and new page table, and then resume the interrupted instruction

- Page Fault Rate $0 \leq p \leq 1$
  - if $p = 0$ no page faults
  - if $p = 1$, every reference is a fault

- Effective Access Time (EAT)

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

# Demand Paging Example

- Memory access time = 200 nanoseconds

- Average page-fault service time = 8 milliseconds

- EAT = (1 – p) x 200 + p (8 milliseconds)

    = (1 – p  x 200 + p x 8,000,000

    = 200 + p x 7,999,800

- If one access out of 1,000 causes a page fault, then

    EAT = 8.2 microseconds.

 This is a slowdown by a factor of 40!!

- If want performance degradation < 10 percent

  - 220 > 200 + 7,999,800 x p
    20 > 7,999,800 x p

  - p < .0000025

  - < one page fault in every 400,000 memory accesses

# Demand Paging Optimizations

- Copy entire process image to swap space at process load time
    - Then page in and out of swap space
    - Used in older BSD Unix

- Demand page in from program binary on disk, but discard rather than paging out when freeing frame
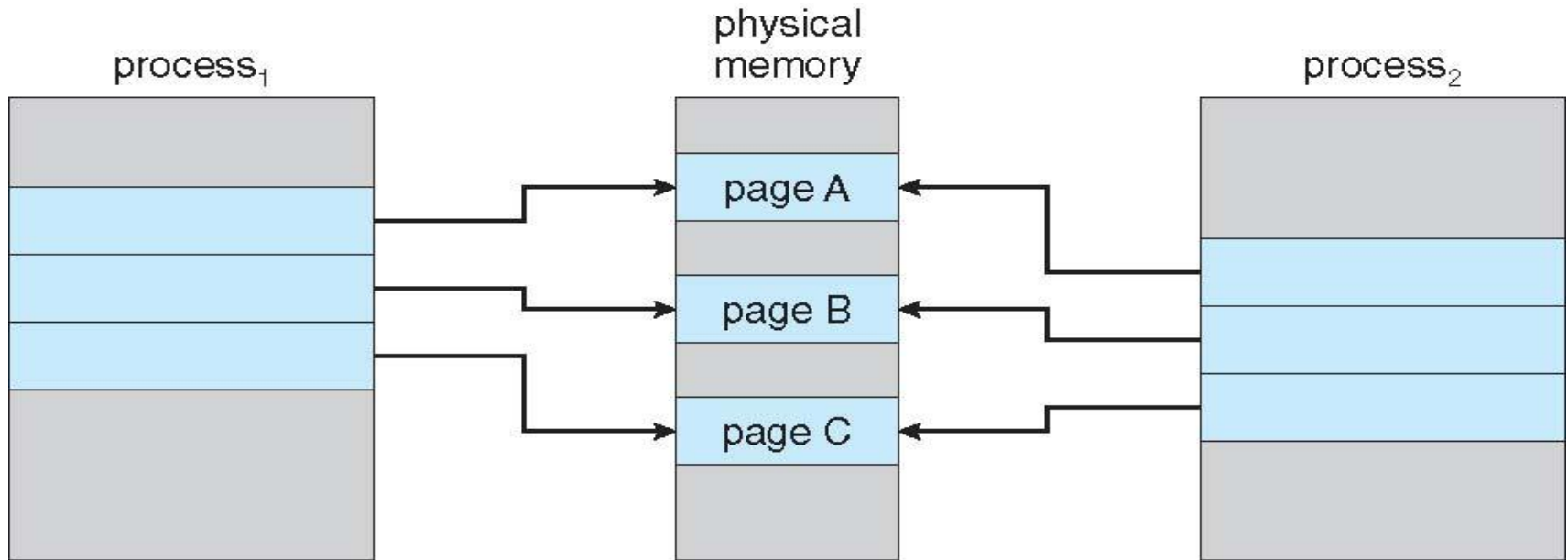    - Used in Solaris and current BSD

# Copy-on-Write

- **Copy-on-Write** (COW) allows both parent and child processes to initially *share* the same pages in memory
  - If either process modifies a shared page, only then is the page copied
- COW allows more efficient process creation as only modified pages are copied
- In general, free pages are allocated from a **pool** of **zero-fill-on-demand** pages
  - Why zero-out a page before allocating it?
- `vfork()` variation on `fork()` system call has parent suspend and child using copy-on-write address space of parent
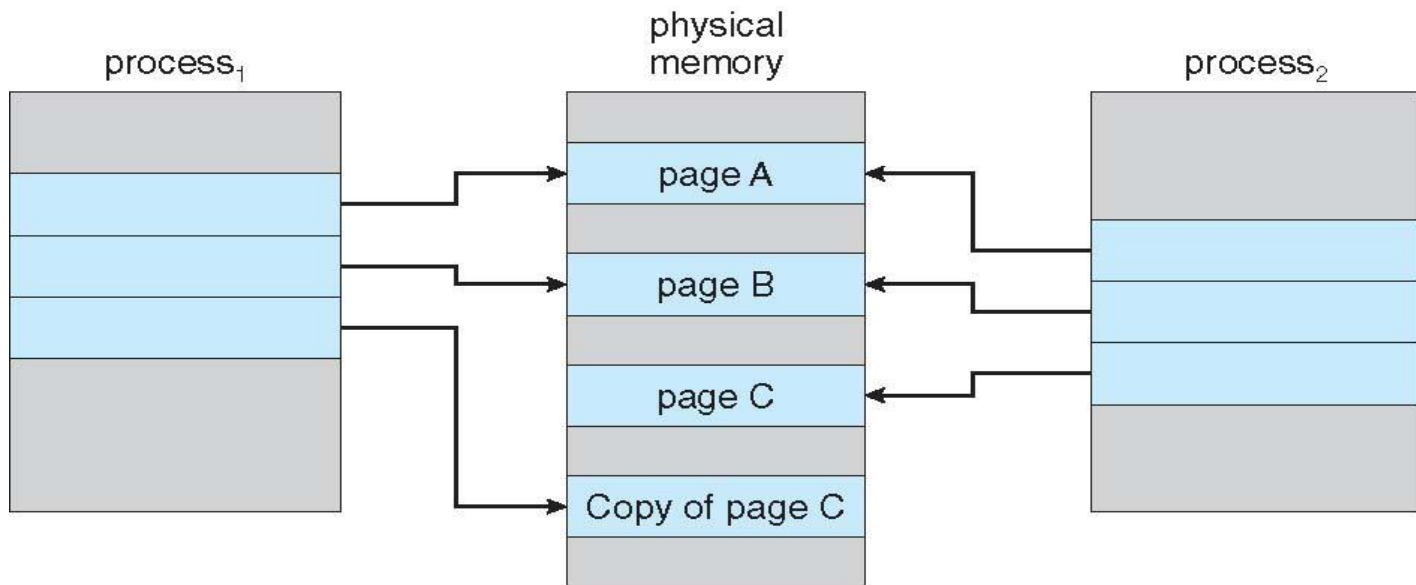  - Designed to have child call `exec()`
  - Very efficient

# Before Process 1 Modifies Page C

# After Process 1 Modifies Page C

# What Happens if There is no Free Frame?

- Used up by process pages

- Also in demand from the kernel, I/O buffers, etc

- How much to allocate to each?

- Page replacement – find some page in memory, but not really in use, page it out
    - Algorithm – terminate? swap out? replace the page?
    - Performance – want an algorithm which will result in minimum number of page faults

- Same page may be brought into memory several times

# Page Replacement
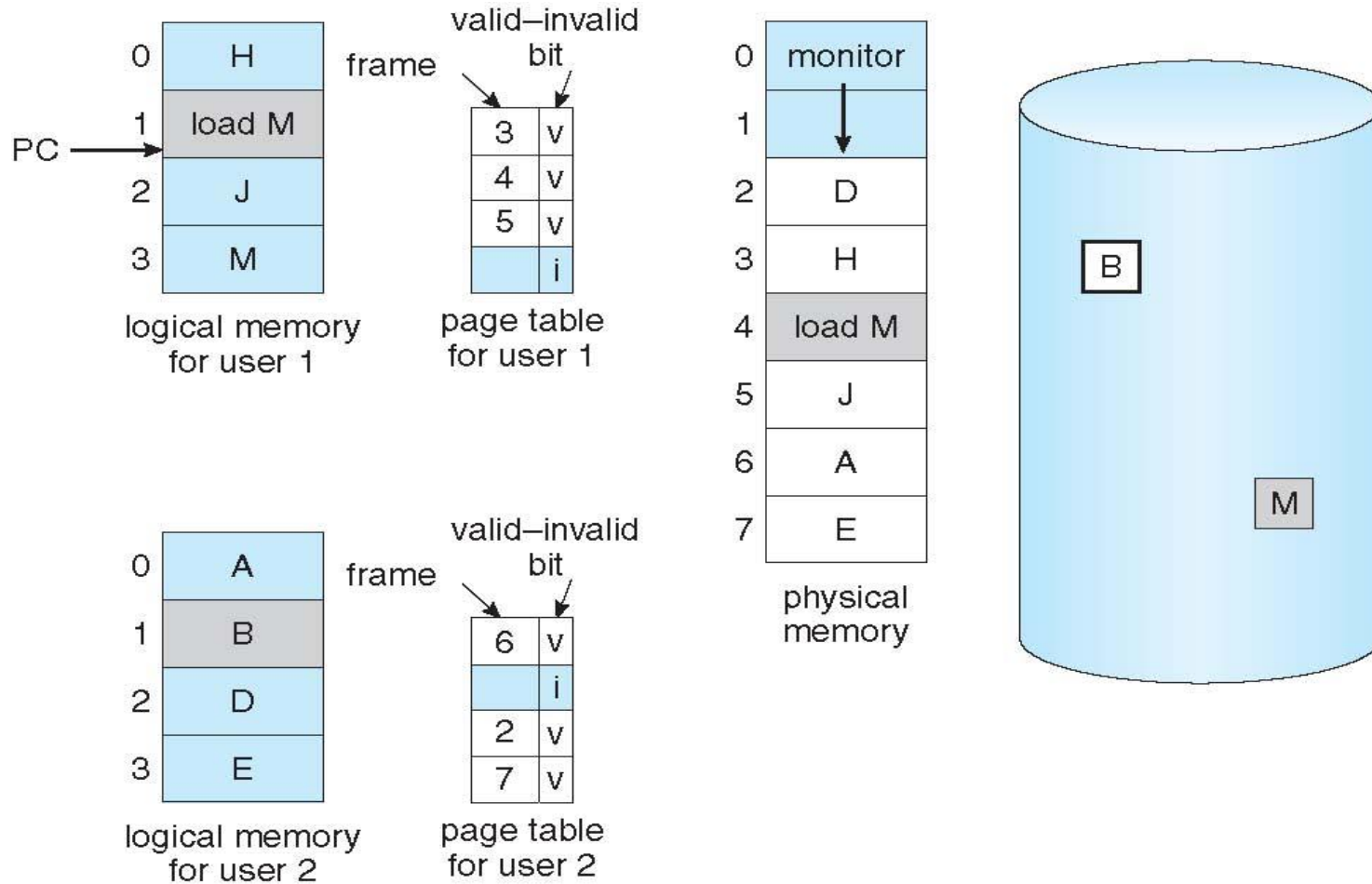
□ Prevent over-allocation of memory by modifying page-fault service routine to include page replacement

□ Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk

□ Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory

# Need For Page Replacement

# Basic Page Replacement

1.   Find the location of the desired page on disk

2.   Find a free frame:
       -  If there is a free frame, use it
       -  If there is no free frame, use a page replacement algorithm to select a **victim frame**
           - Write victim frame to disk if dirty

3.   Bring  the desired page into the (newly) free frame; update the page and frame tables

4.   Continue the process by restarting the instruction that caused the trap

Note now potentially 2 page transfers for page fault – increasing EAT
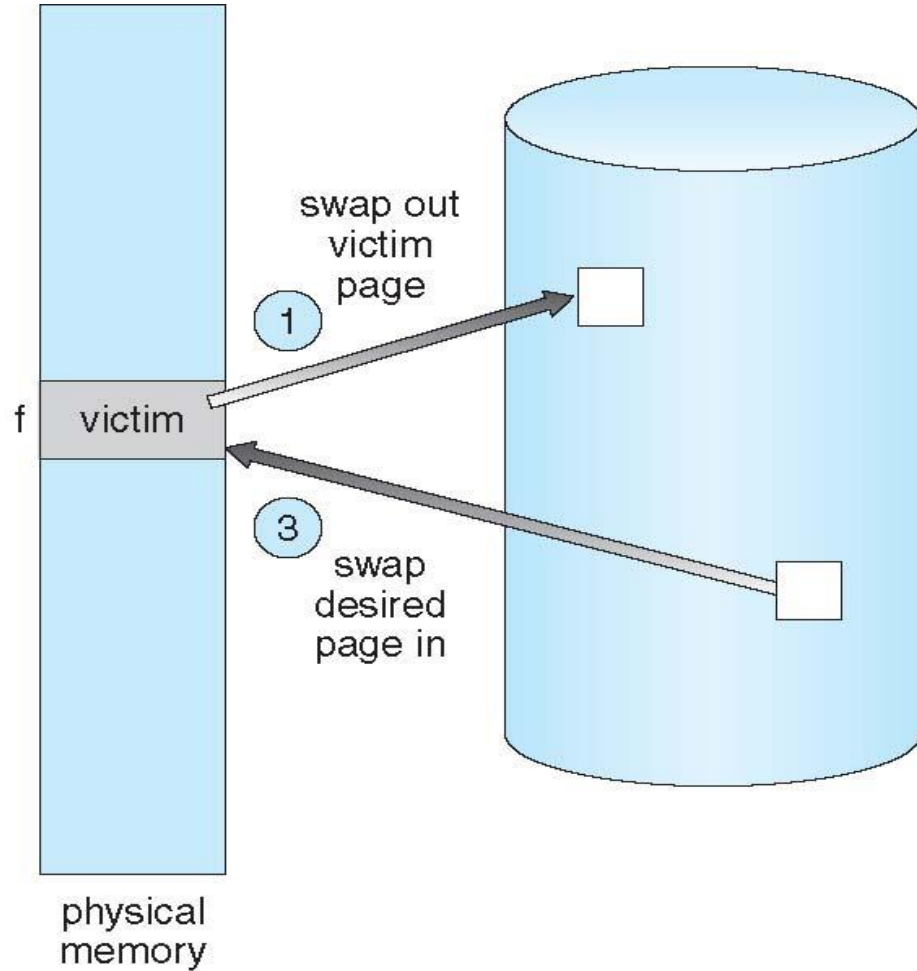
# Page Replacement

# Page and Frame Replacement Algorithms

- **Frame-allocation algorithm** determines
  - How many frames to give each process
  - Which frames to replace
- **Page-replacement algorithm**
  - Want lowest page-fault rate on both first access and re-access

- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
  - String is just page numbers, not full addresses
  - Repeated access to the same page does not cause a page fault
- In all our examples, the reference string is

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

# Graph of Page Faults Versus The Number of Frames

# First-In-First-Out (FIFO) Algorithm

- Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**

- 3 frames (3 pages can be in memory at a time per process)

| 1 | 7 | 2 | 4 0 7 |
|---|---|---|-------|
| 2 | 0 | 3 | 2 1 0 |
| 3 | 1 | 0 | 3 2 1 |

15 page faults

- Can vary by reference string: consider 1,2,3,4,1,2,5,1,2,3,4,5

  - Adding more frames can cause more page faults!

    ‣ **Belady's Anomaly**

- How to track ages of pages?

  - Just use a FIFO queue

# FIFO Page Replacement



reference string
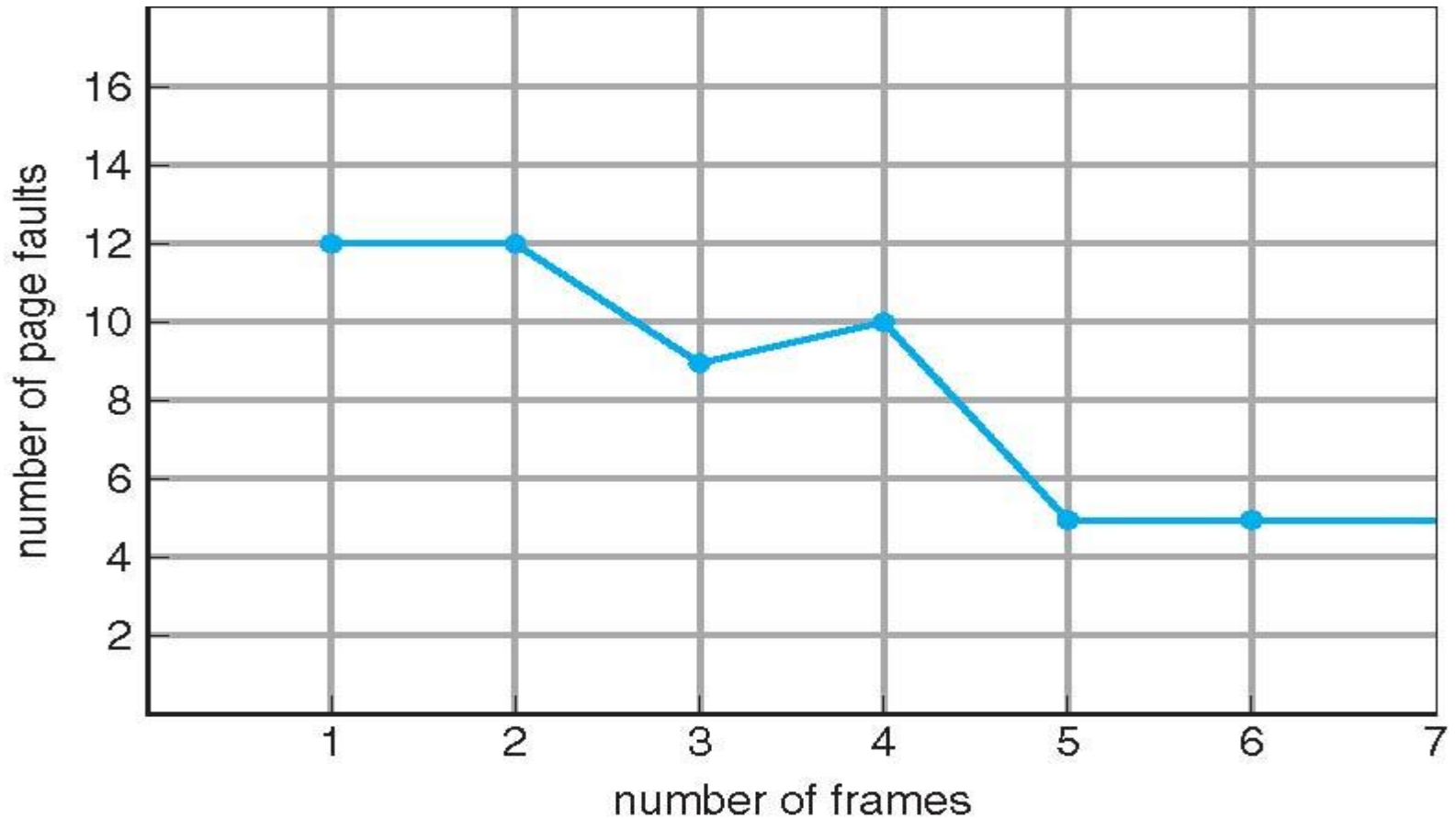
7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1

page frames

# FIFO Illustrating Belady's Anomaly

# Optimal Algorithm

- Replace page that will not be used for longest period of time
    - 9 is optimal for the example on the next slide

- How do you know this?
    - Can't read the future

- Used for measuring how well your algorithm performs

# Optimal Page Replacement

# Least Recently Used (LRU) Algorithm

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

reference string

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



page frames

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

# LRU Algorithm (Cont.)

- Counter implementation
    - Every page entry has a counter; every time page is referenced through this entry, copy the clock into the counter
    - When a page needs to be changed, look at the counters to find smallest value
        - ▸ Search through table needed
- Stack implementation
    - Keep a stack of page numbers in a double link form:
    - Page referenced:
        - ▸ move it to the top
        - ▸ requires 6 pointers to be changed
    - But each update more expensive
    - No search for replacement
- LRU and OPT are cases of **stack algorithms** that don't have Belady's Anomaly

# Use Of A Stack to Record The Most Recent Page References



reference string

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

stack before a

| 2 |
| 1 |
| 0 |
| 7 |
| 4 |

stack after b

| 7 |
| 2 |
| 1 |
| 0 |
| 4 |

# LRU Approximation Algorithms

- LRU needs special hardware and still slow
- **Reference bit**
    - With each page associate a bit, initially = 0
    - When page is referenced bit set to 1
    - Replace any with reference bit = 0 (if one exists)
        - We do not know the order, however
- **Second-chance algorithm**
    - Generally FIFO, plus hardware-provided reference bit
    - Clock replacement
    - If page to be replaced has
        - Reference bit = 0 -> replace it
        - reference bit = 1 then:
            - set reference bit 0, leave page in memory
            - replace next page, subject to same rules

# Second-Chance (clock) Page-Replacement Algorithm

# Counting Algorithms

- Keep a counter of the number of references that have been made to each page

  - Not common

- **LFU Algorithm**:  replaces page with smallest count

- **MFU Algorithm**: based on the argument that the page with the smallest count was probably just brought in and has yet to be used

# Page-Buffering Algorithms

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

# Applications and Page Replacement

- All of these algorithms have OS guessing about future page access
- Some applications have better knowledge – i.e. databases
- Memory intensive applications can cause double buffering
  - OS keeps copy of page in memory as I/O buffer
  - Application keeps page in memory for its own work
- Operating system can given direct access to the disk, getting out of the way of the applications
  - **Raw disk** mode
- Bypasses buffering, locking, etc

# Allocation of Frames

- Each process needs *minimum* number of frames
- Example: IBM 370 – 6 pages to handle SS MOVE instruction:
    - instruction is 6 bytes, might span 2 pages
    - 2 pages to handle *from*
    - 2 pages to handle *to*
- *Maximum* of course is total frames in the system
- Two major allocation schemes
    - fixed allocation
    - priority allocation
- Many variations

# Fixed Allocation

☐ Equal allocation – For example, if there are 100 frames (after allocating frames for the OS) and 5 processes, give each process 20 frames

    ☐ Keep some as free frame buffer pool

☐ Proportional allocation – Allocate according to the size of process

    ☐ Dynamic as degree of multiprogramming, process sizes change

$-\ s_i = $ size of process $p_i$

$-\ S = \sum s_i$

$-\ m = $ total number of frames

$-\ a_i = $ allocation for $p_i = \dfrac{s_i}{S} \times m$

$m = 64$

$s = 10$

$s_2 = 127$

$a_1 = \dfrac{10}{137} \times 64 \approx 5$

$a_2 = \dfrac{127}{137} \times 64 \approx 59$

# **Priority Allocation**

- Use a proportional allocation scheme using priorities rather than size

- If process $P_i$ generates a page fault,
    - select for replacement one of its frames
    - select for replacement a frame from a process with lower priority number

# Global vs. Local Allocation

□ **Global replacement** – process selects a replacement frame from the set of all frames; one process can take a frame from another

  □ But then process execution time can vary greatly

  □ But greater throughput so more common

□ **Local replacement** – each process selects from only its own set of allocated frames

  □ More consistent per-process performance

  □ But possibly underutilized memory

# Non-Uniform Memory Access

- So far all memory accessed equally

- Many systems are NUMA – speed of access to memory varies
    - Consider system boards containing CPUs and memory, interconnected over a system bus

- Optimal performance comes from allocating memory "close to" the CPU on which the thread is scheduled
    - And modifying the scheduler to schedule the thread on the same system board when possible
    - Solved by Solaris by creating **lgroups**
        - ▸ Structure to track CPU / Memory low latency groups
        - ▸ Used my schedule and pager
        - ▸ When possible schedule all threads of a process and allocate all memory for that process within the lgroup
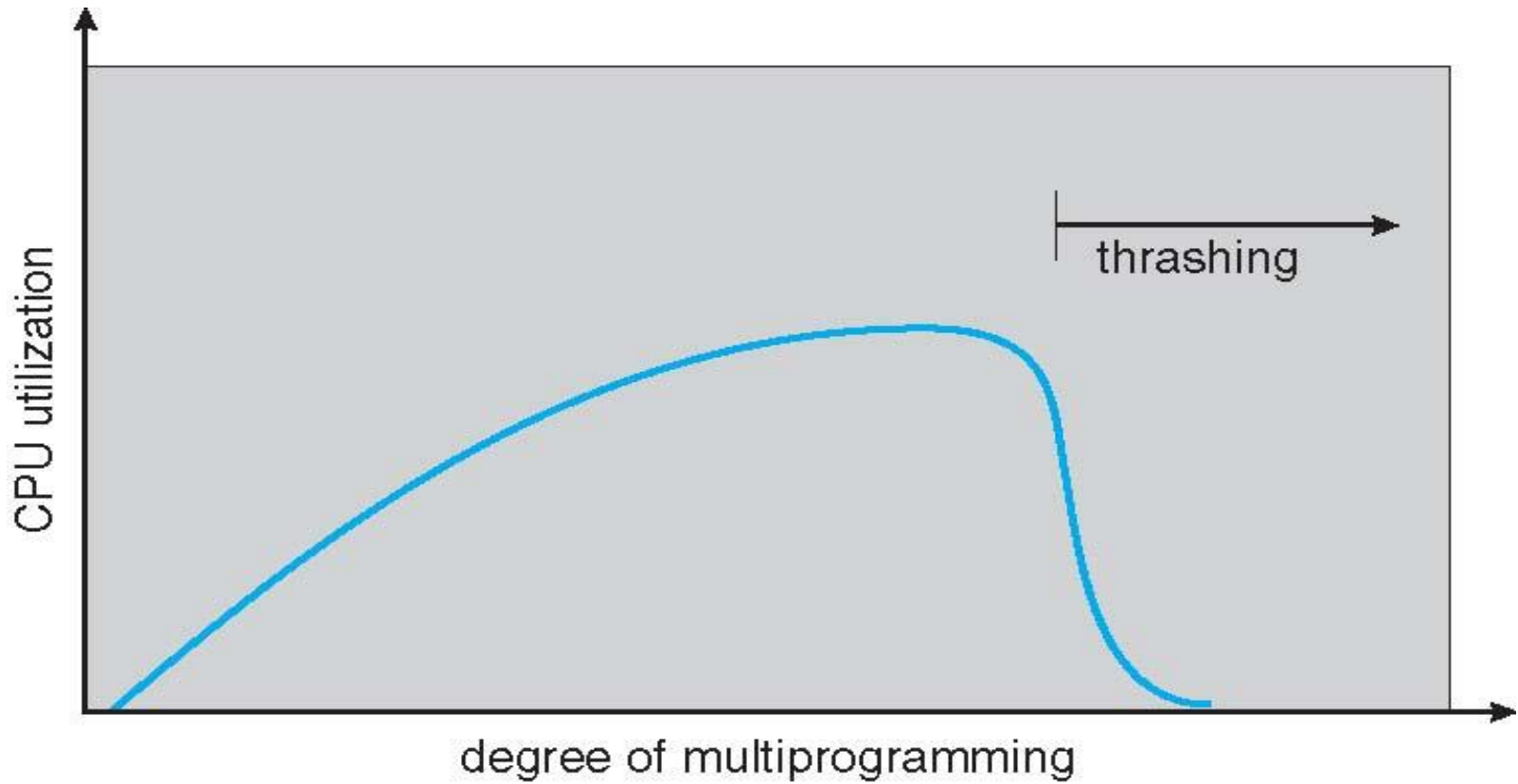
# Thrashing

- If a process does not have "enough" pages, the page-fault rate is very high

    - Page fault to get page

    - Replace existing frame

    - But quickly need replaced frame back

    - This leads to:

        ▸ Low CPU utilization

        ▸ Operating system thinking that it needs to increase the degree of multiprogramming

        ▸ Another process added to the system

- **Thrashing** ≡ a process is busy swapping pages in and out

# Thrashing (Cont.)

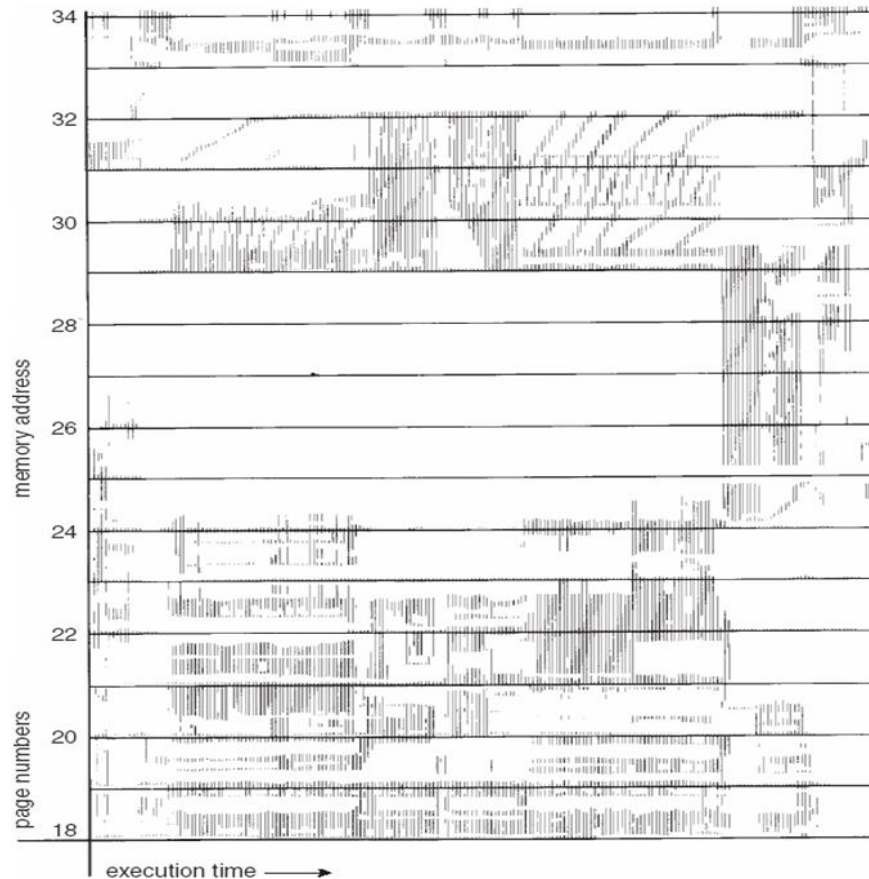# Demand Paging and Thrashing

□ Why does demand paging work?
**Locality model**

    □ Process migrates from one locality to another

    □ Localities may overlap


□ Why does thrashing occur?
$\Sigma$ size of locality > total memory size

    □ Limit effects by using local or priority page replacement

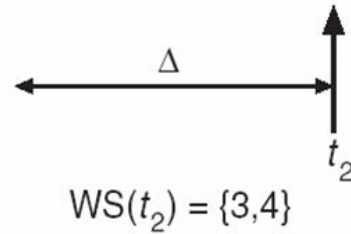# Locality In A Memory-Reference Pattern

# Working-Set Model

- $\Delta \equiv$ working-set window $\equiv$ a fixed number of page references
  Example: 10,000 instructions

- $WSS_i$ (working set of Process $P_i$) =
  total number of pages referenced in the most recent $\Delta$ (varies in time)

  - if $\Delta$ too small will not encompass entire locality

  - if $\Delta$ too large will encompass several localities

  - if $\Delta = \infty \Rightarrow$ will encompass entire program

- $D = \Sigma\ WSS_i \equiv$ total demand frames

  - Approximation of locality

- if $D > m \Rightarrow$ Thrashing

- Policy if $D > m$, then suspend or swap out one of the processes

# Working-set model

page reference table

... 2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 ...

$\Delta$          $t_1$

$\Delta$          $t_2$

$WS(t_1) = \{1,2,5,6,7\}$          $WS(t_2) = \{3,4\}$

# Keeping Track of the Working Set

- Approximate with interval timer + a reference bit

- Example: $\Delta$ = 10,000
  - Timer interrupts after every 5000 time units
  - Keep in memory 2 bits for each page
  - Whenever a timer interrupts copy and sets the values of all reference bits to 0
  - If one of the bits in memory = 1 $\Rightarrow$ page in working set

- Why is this not completely accurate?

- Improvement = 10 bits and interrupt every 1000 time units
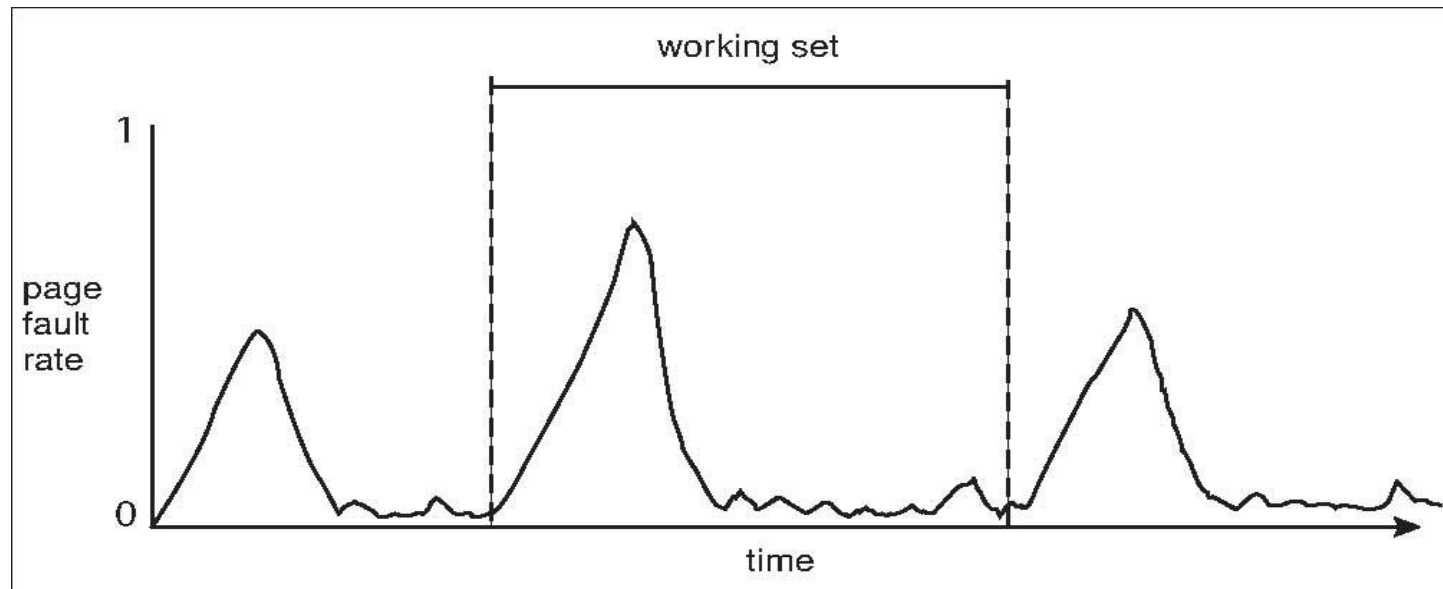
# Page-Fault Frequency

- More direct approach than WSS

- Establish "acceptable" **page-fault frequency** rate and use local replacement policy

  - If actual rate too low, process loses frame

  - If actual rate too high, process gains frame

# Working Sets and Page Fault Rates

# Memory-Mapped Files

- Memory-mapped file I/O allows file I/O to be treated as routine memory access by **mapping** a disk block to a page in memory

- A file is initially read using demand paging
  - A page-sized portion of the file is read from the file system into a physical page
  - Subsequent reads/writes to/from the file are treated as ordinary memory accesses

- Simplifies and speeds file access by driving file I/O through memory rather than `read()` and `write()` system calls

- Also allows several processes to map the same file allowing the pages in memory to be shared

- But when does written data make it to disk?
  - Periodically and / or at file `close()` time
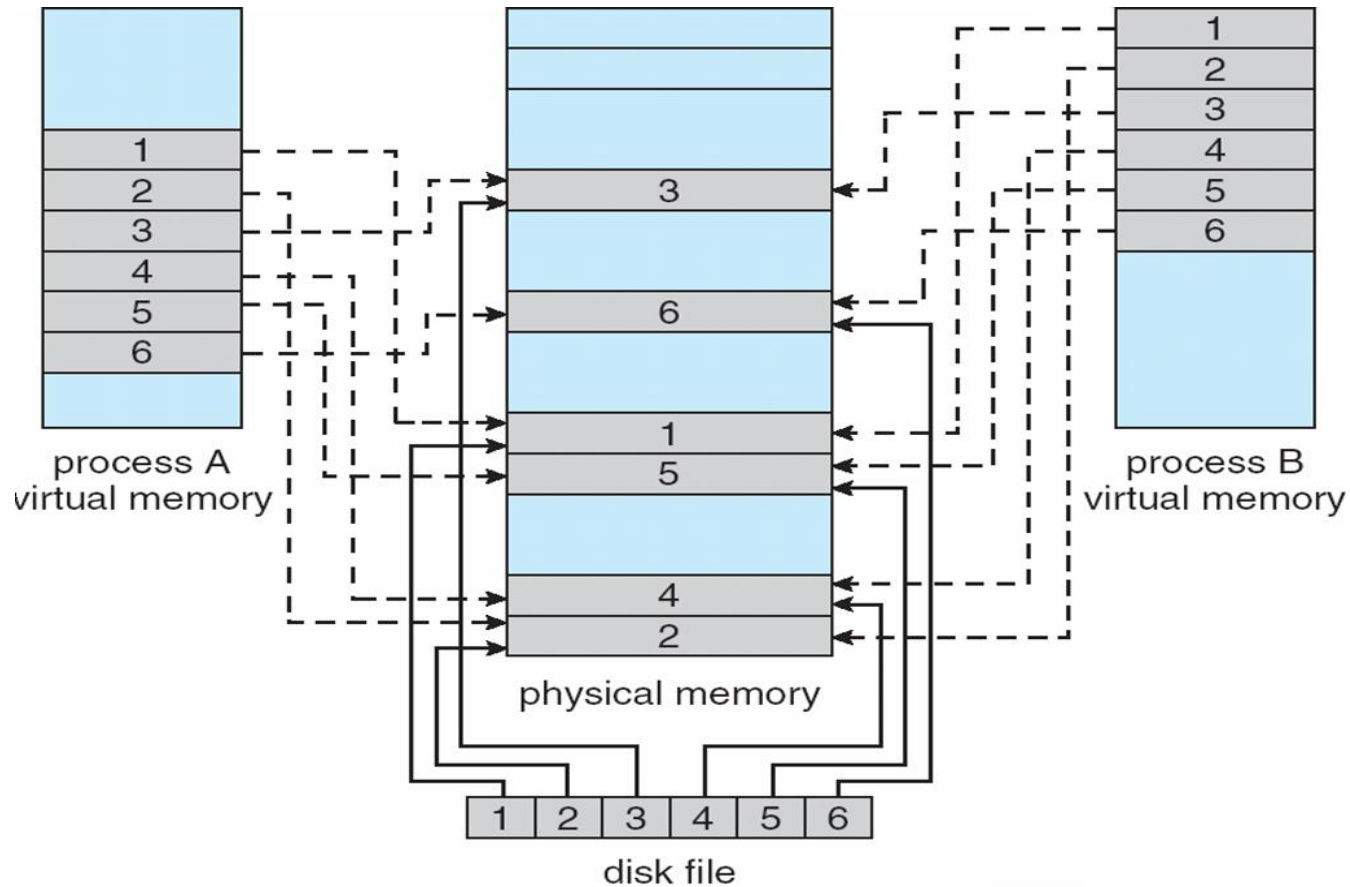  - For example, when the pager scans for dirty pages

# Memory-Mapped File Technique for all I/O

- Some OSes uses memory mapped files for standard I/O
- Process can explicitly request memory mapping a file via `mmap()` system call
  - Now file mapped into process address space
- For standard I/O (`open()`, `read()`, `write()`, `close()`), mmap anyway
  - But map file into kernel address space
  - Process still does read() and write()
    - ▸ Copies data to and from kernel space and user space
  - Uses efficient memory management subsystem
    - ▸ Avoids needing separate subsystem
- COW can be used for read/write non-shared pages
- Memory mapped files can be used for shared memory (although again via separate system calls)

# Memory Mapped Files

# Memory-Mapped Shared Memory in Windows

# File-System Interface

- File Concept
- Access Methods
- Disk and Directory Structure
- File-System Mounting
- File Sharing
- Protection

# File Concept

- Contiguous logical address space
- Types:
    - Data
        - numeric
        - character
        - binary
    - Program
- Contents defined by file's creator
    - Many types
        - Consider **text file, source file, executable file**

# File Attributes

- **Name** – only information kept in human-readable form
- **Identifier** – unique tag (number) identifies file within file system
- **Type** – needed for systems that support different types
- **Location** – pointer to file location on device
- **Size** – current file size
- **Protection** – controls who can do reading, writing, executing
- **Time, date, and user identification** – data for protection, security, and usage monitoring
- Information about files are kept in the directory structure, which is maintained on the disk
- Many variations, including extended file attributes such as file checksum
- Information kept in the directory structure

# File info Window on Mac OS X

# File Operations

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

# Open Files

- Several pieces of data are needed to manage open files:

  - **Open-file table**: tracks open files

  - File pointer:  pointer to last read/write location, per process that has the file open

  - **File-open count**: counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it

  - Disk location of the file: cache of data access information

  - Access rights: per-process access mode information

# Open File Locking

- Provided by some operating systems and file systems

  - Similar to reader-writer locks

  - **Shared lock** similar to reader lock – several processes can acquire concurrently

  - **Exclusive lock** similar to writer lock

- Mediates access to a file

- Mandatory or advisory:

  - **Mandatory** – access is denied depending on locks held and requested

  - **Advisory** – processes can find status of locks and decide what to do

# File Locking Example – Java API

```java
import java.io.*;
import java.nio.channels.*;
public class LockingExample {
    public static final boolean EXCLUSIVE = false;
    public static final boolean SHARED = true;
    public static void main(String arsg[]) throws IOException {
        FileLock sharedLock = null;
        FileLock exclusiveLock = null;
        try {
                RandomAccessFile raf = new RandomAccessFile("file.txt", "rw");
                // get the channel for the file
                FileChannel ch = raf.getChannel();
                // this locks the first half of the file - exclusive
                exclusiveLock = ch.lock(0, raf.length()/2, EXCLUSIVE);
                /** Now modify the data . . . */
                // release the lock
                exclusiveLock.release();
```

# File Locking Example – Java API (Cont.)

```
                    // this locks the second half of the file - shared
                    sharedLock = ch.lock(raf.length()/2+1, raf.length(),
                                        SHARED);
                    /** Now read the data . . . */
                    // release the lock
                    sharedLock.release();
            } catch (java.io.IOException ioe) {
                    System.err.println(ioe);
            }finally {

                    if (exclusiveLock != null)
                    exclusiveLock.release();
                    if (sharedLock != null)
                    sharedLock.release();
            }
        }
    }
```

# File Types – Name, Extension

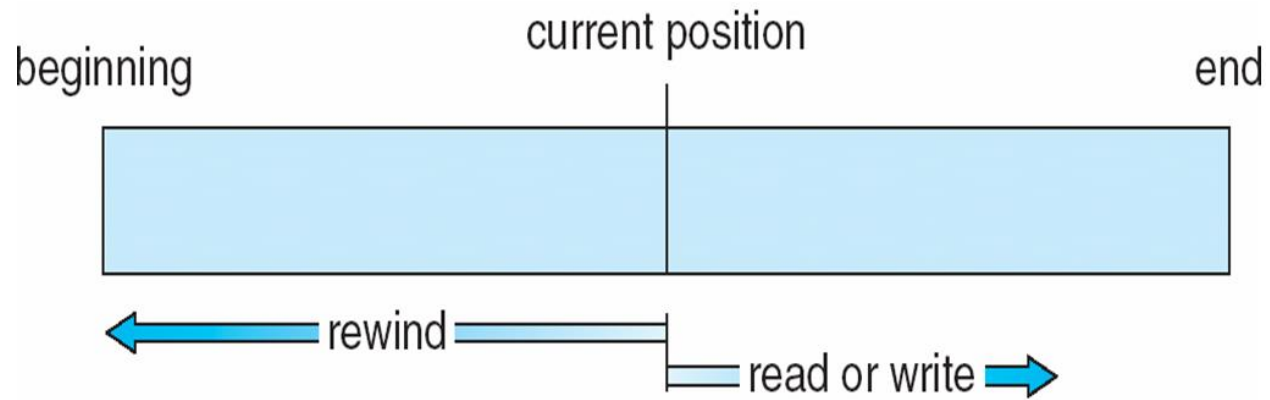| file type | usual extension | function |
|---|---|---|
| executable | exe, com, bin or none | ready-to-run machine-language program |
| object | obj, o | compiled, machine language, not linked |
| source code | c, cc, java, pas, asm, a | source code in various languages |
| batch | bat, sh | commands 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 com-pressed, for archiving or storage |
| multimedia | mpeg, mov, rm, mp3, avi | binary file containing audio or A/V information |

# File Structure

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

# Sequential-access File

# Access Methods

☐ **Sequential Access**

```
read next
write next
reset
```
no read after last write
(rewrite)

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

```
read n
write n
position to n
        read next
        write next
rewrite n
```

*n* = relative block number


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

  ☐ See allocation problem in Ch 12

# Simulation of Sequential Access on Direct-access File

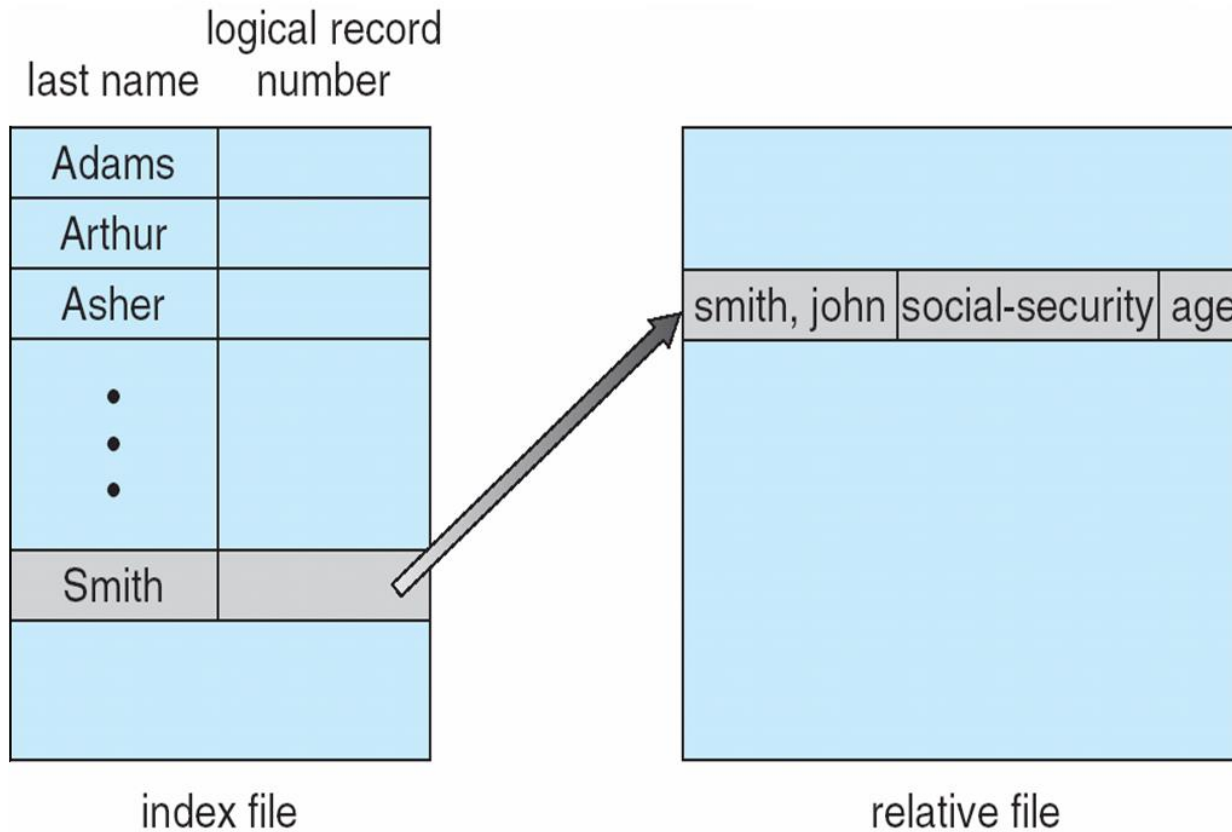| sequential access | implementation for direct access |
|---|---|
| reset | $cp = 0;$ |
| read next | read $cp$; <br> $cp = cp + 1;$ |
| write next | write $cp$; <br> $cp = cp + 1;$ |

# Other Access Methods

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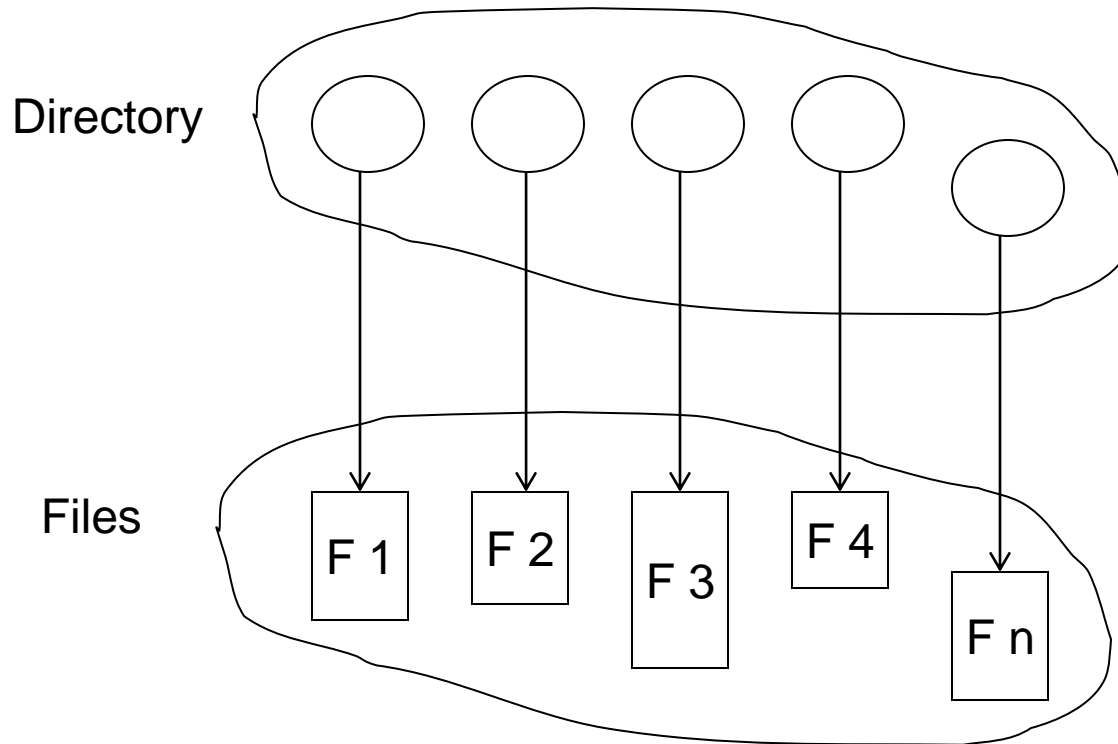- VMS operating system provides index and relative files as another example (see next slide)

# Example of Index and Relative Files



index file

relative file

# Directory Structure

☐ A collection of nodes containing information about all files



Directory

Files
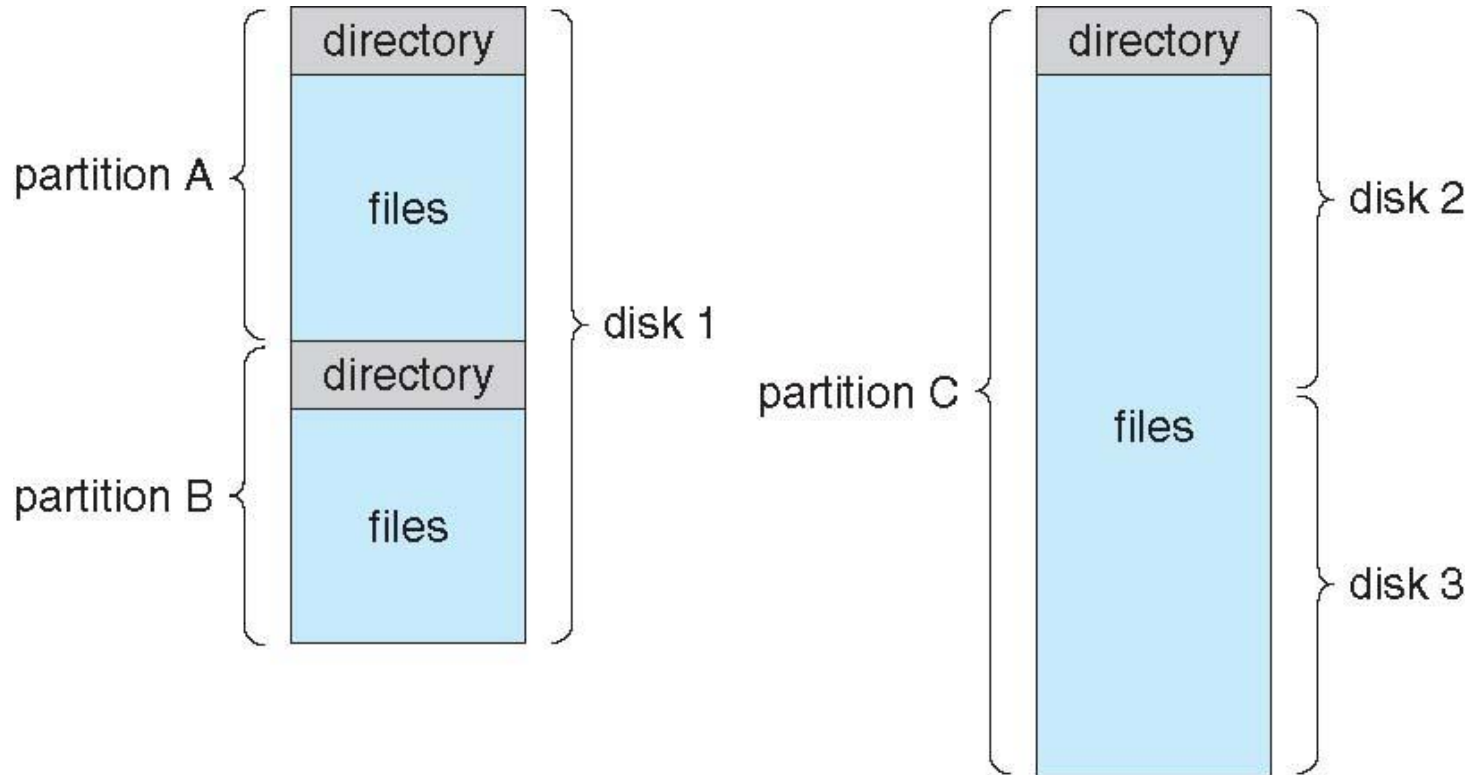
Both the directory structure and the files reside on disk

# Disk Structure

- Disk can be subdivided into **partitions**

- Disks or partitions can be **RAID** protected against failure

- Disk or partition can be used **raw** – without a file system, or **formatted** with a file system

- Partitions also known as minidisks, slices

- Entity containing file system known as a **volume**

- Each volume containing file system also tracks that file system's info in **device directory** or **volume table of contents**

- As well as **general-purpose file systems** there are many **special-purpose file systems**, frequently all within the same operating system or computer
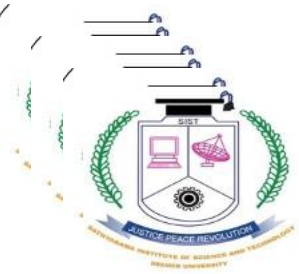
# A Typical File-system Organization

# Types of File Systems

- We mostly talk of general-purpose file systems

- But systems frequently have may file systems, some general- and some special- purpose

- Consider Solaris has

    - tmpfs – memory-based volatile FS for fast, temporary I/O

    - objfs – interface into kernel memory to get kernel symbols for debugging

    - ctfs – contract file system for managing daemons

    - lofs – loopback file system allows one FS to be accessed in place of another

    - procfs – kernel interface to process structures

    - ufs, zfs – general purpose file systems

# Operations Performed on Directory

- Search for a file

- Create a file

- Delete a file

- List a directory

- Rename a file

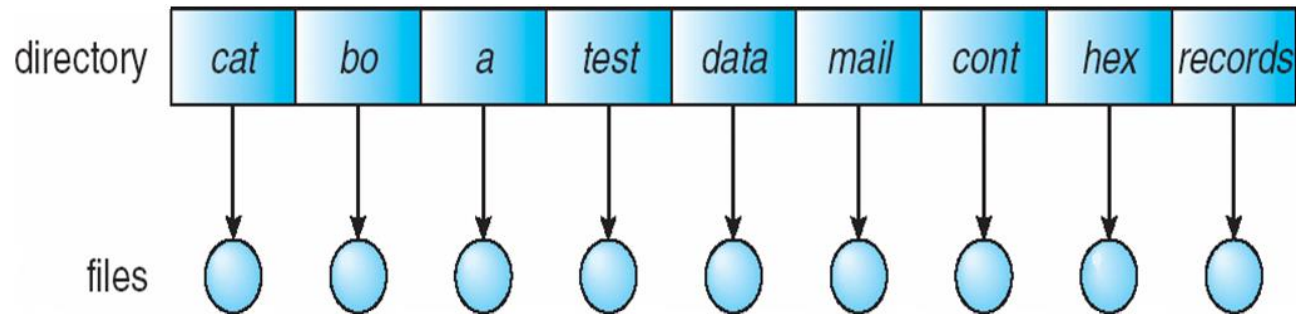- Traverse the file system

# Directory Organization

rectory is organized logically  to obtain

- Efficiency – locating a file quickly
- Naming – convenient to users
    - Two users can have same name for different files
    - The same file can have several different names
- Grouping – logical grouping of files by properties, (e.g., all Java programs, all games, …)

# Single-Level Directory

- A single directory for all users



- Naming problem
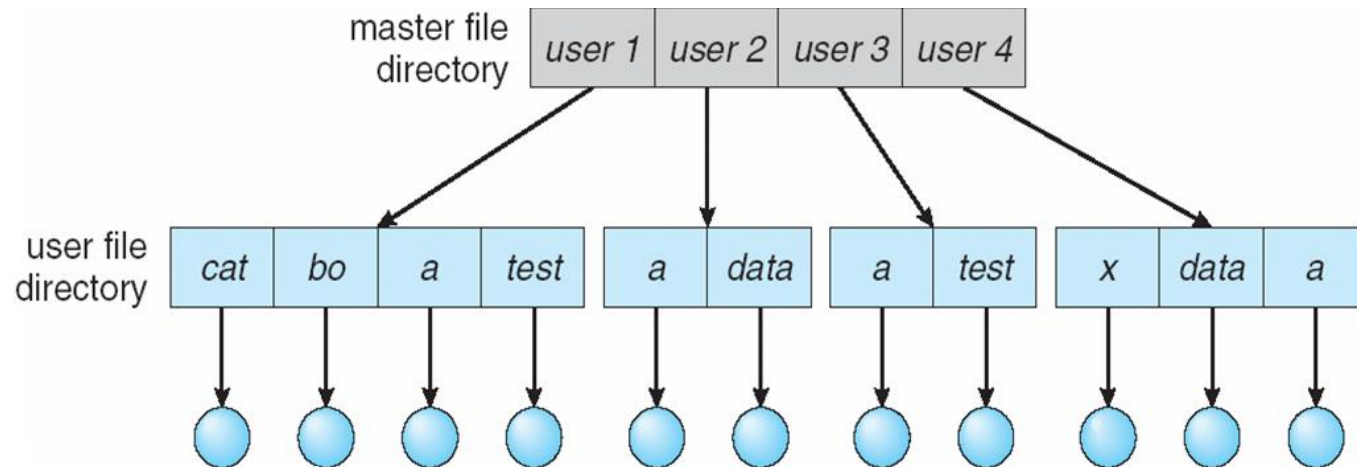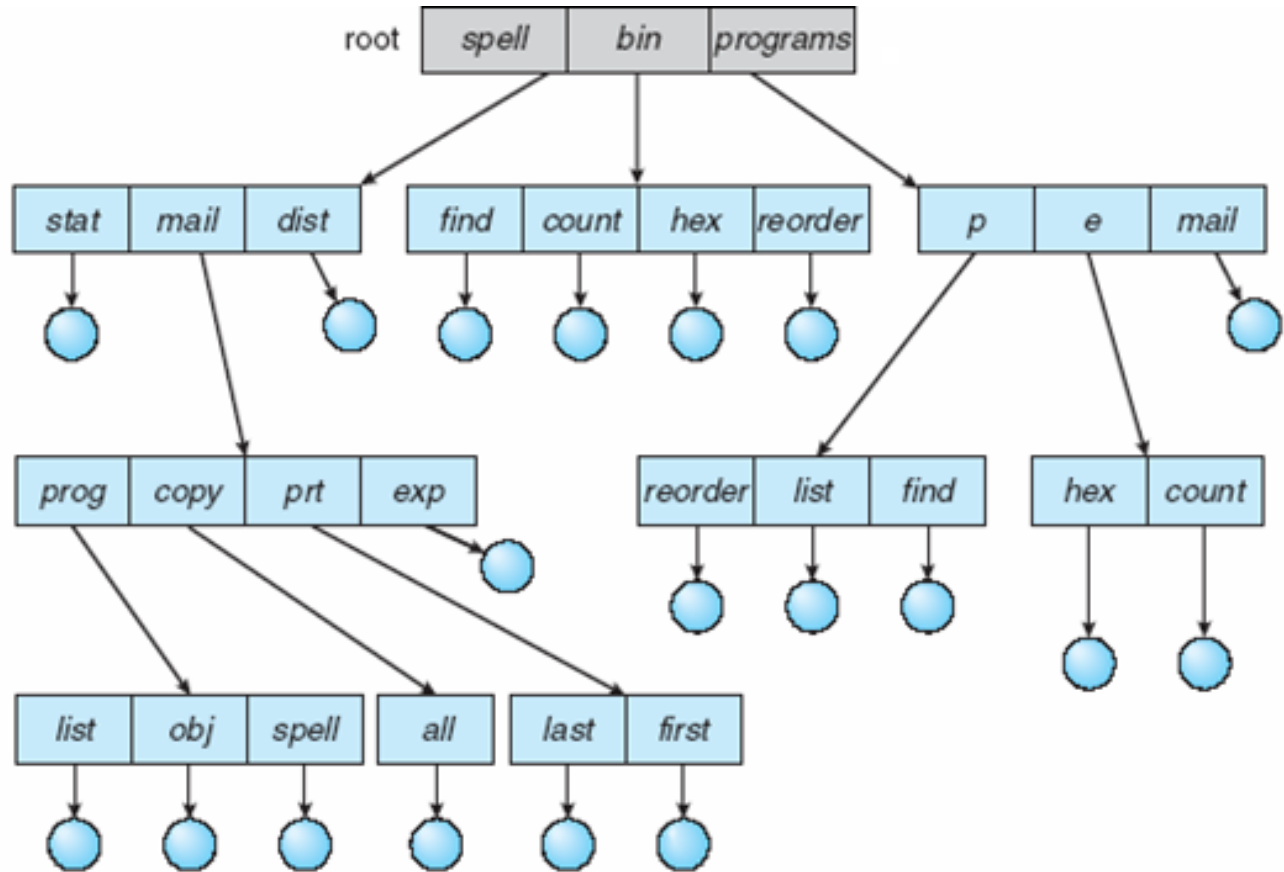- Grouping problem

# Two-Level Directory

☐ Separate directory for each user



☐ Path name

☐ Can have the same file name for different user

☐ Efficient searching

☐ No grouping capability

# Tree-Structured Directories

# Tree-Structured Directories (Cont.)

- Efficient searching

- Grouping Capability

- Current directory (working directory)
    - `cd /spell/mail/prog`
    - `type list`

# Tree-Structured Directories (Cont)

- **Absolute** or **relative** path name
- Creating a new file is done in current directory
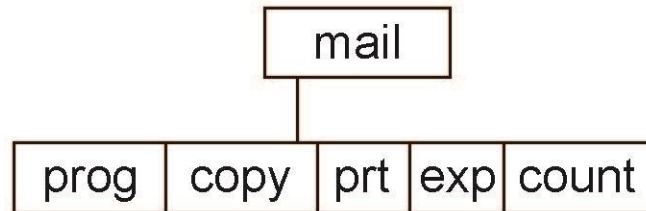- Delete a file

```
rm <file-name>
```

- Creating a new subdirectory is done in current directory

```
mkdir <dir-name>
```

Example:  if in current directory  `/mail`

```
mkdir count
```



Deleting "mail" $\Rightarrow$ deleting the entire subtree rooted by "mail"

# Acyclic-Graph Directories

- Have shared subdirectories and files

# Acyclic-Graph Directories (Cont.)

- Two different names (aliasing)
- If **dict** deletes **list** $\Rightarrow$ dangling pointer

  Solutions:

  - Backpointers, so we can delete all pointers
    Variable size records a problem
  - Backpointers using a daisy chain organization
  - Entry-hold-count solution

- New directory entry type

  - **Link** – another name (pointer) to an existing file
  - **Resolve the link** – follow pointer to locate the file

# General Graph Directory

# General Graph Directory (Cont.)

- How do we guarantee no cycles?
  - Allow only links to file not subdirectories
  - **Garbage collection**
  - Every time a new link is added use a cycle detection algorithm to determine whether it is OK

# File Sharing

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

# File Sharing – Remote File Systems

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

# File Sharing – Failure Modes

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

# File Sharing – Consistency Semantics

- Specify how multiple users are to access a shared file simultaneously

  - Similar to Ch 5 process synchronization algorithms

    - Tend to be less complex due to disk I/O and network latency (for remote file systems

  - Andrew File System (AFS) implemented complex remote file sharing semantics

  - Unix file system (UFS) implements:

    - Writes to an open file visible immediately to other users of the same open file

    - Sharing file pointer to allow multiple users to read and write concurrently

  - AFS has session semantics

    - Writes only visible to sessions starting after the file is closed

# Protection

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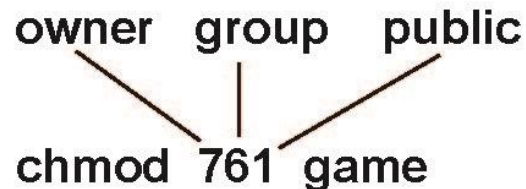    - **Delete**
    - **List**

# Access Lists and Groups

- Mode of access: read, write, execute
- Three classes of users on Unix / Linux

|  |  |  | RWX |
|---|---|---|---|
| a) **owner access** | 7 | ⇒ | 1 1 1 |
| b) **group access** | 6 | ⇒ | RWX<br>1 1 0 |
| c) **public access** | 1 | ⇒ | RWX<br>0 0 1 |

- Ask manager to create a group (unique name), say G, and add some users to the group.
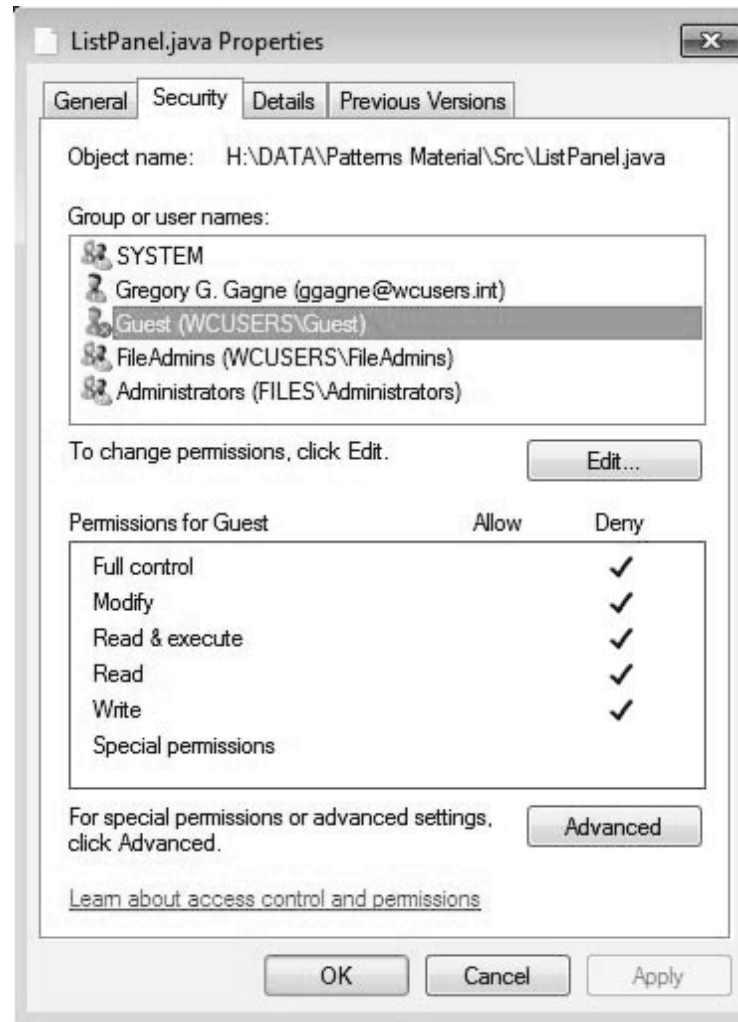- For a particular file (say *game*) or subdirectory, define an appropriate access.

owner   group   public

chmod  761  game

Attach a group to a file

```
chgrp      G      game
```

# Windows 7 Access-Control List Management

# A Sample UNIX Directory Listing

```
-rw-rw-r--     1 pbg   staff     31200   Sep 3 08:30    intro.ps
drwx------     5 pbg   staff       512   Jul 8 09.33    private/
drwxrwxr-x     2 pbg   staff       512   Jul 8 09:35    doc/
drwxrwx---     2 pbg   student     512   Aug 3 14:13    student-proj/
-rw-r--r--     1 pbg   staff      9423   Feb 24 2003    program.c
-rwxr-xr-x     1 pbg   staff     20471   Feb 24 2003    program
drwx--x--x     4 pbg   faculty     512   Jul 31 10:31   lib/
drwx------     3 pbg   staff      1024   Aug 29 06:52   mail/
drwxrwxrwx     3 pbg   staff       512   Jul 8 09:35    test/
```

# File System Implementation
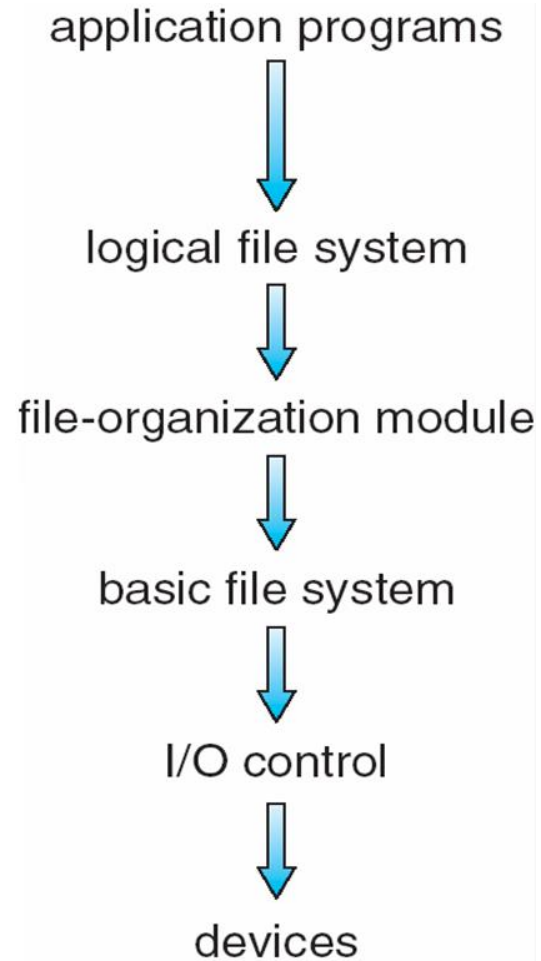
- File-System Structure
- File-System Implementation

# File-System Structure

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

# Layered File System

# File System Layers

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

# File System Layers (Cont.)

- **Logical file system** manages metadata information

  - Translates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)

  - Directory management

  - Protection

- Layering useful for reducing complexity and redundancy, but adds overhead and can decrease performanceTranslates file name into file number, file handle, location by maintaining file control blocks (**inodes** in UNIX)

  - Logical layers can be implemented by any coding method according to OS designer

# File System Layers (Cont.)

- Many file systems, sometimes many within an operating system
    - Each with its own format (CD-ROM is ISO 9660; Unix has **UFS**, FFS; Windows has FAT, FAT32, NTFS as well as floppy, CD, DVD Blu-ray, Linux has more than 40 types, with **extended file system** ext2 and ext3 leading; plus distributed file systems, etc.)
    - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE

# File-System Implementation

- We have system calls at the API level, but how do we implement their functions?

    - On-disk and in-memory structures

- **Boot control block** contains info needed by system to boot OS from that volume

    - Needed if volume contains OS, usually first block of volume

- **Volume control block** (**superblock, master file table**) contains volume details

    - Total # of blocks, # of free blocks, block size, free block pointers or array

- Directory structure organizes the files

    - Names and inode numbers, master file table

# File-System Implementation (Cont.)

- Per-file **File Control Block** (**FCB**) contains many details about the file
  - inode number, permissions, size, dates
  - NFTS stores into in master file table using relational DB structures

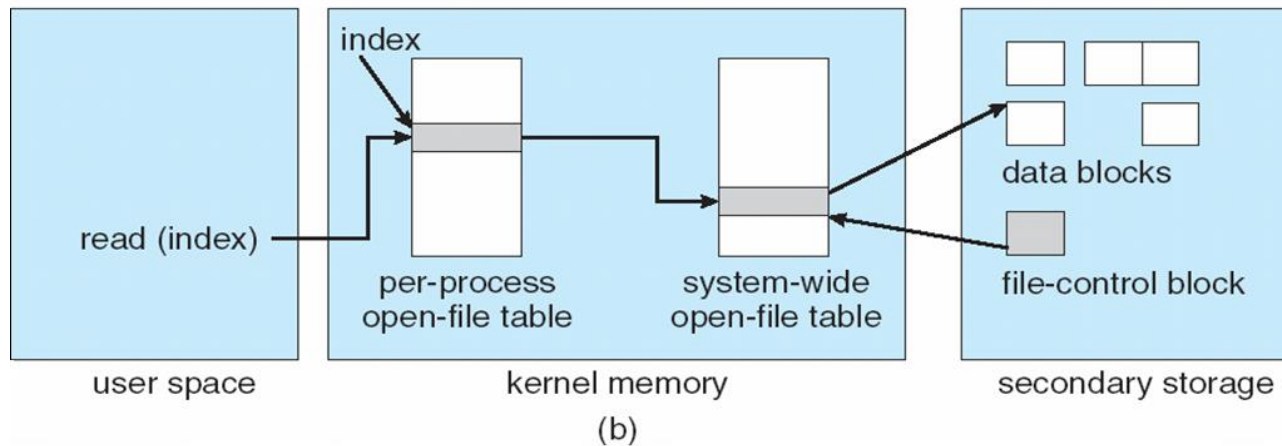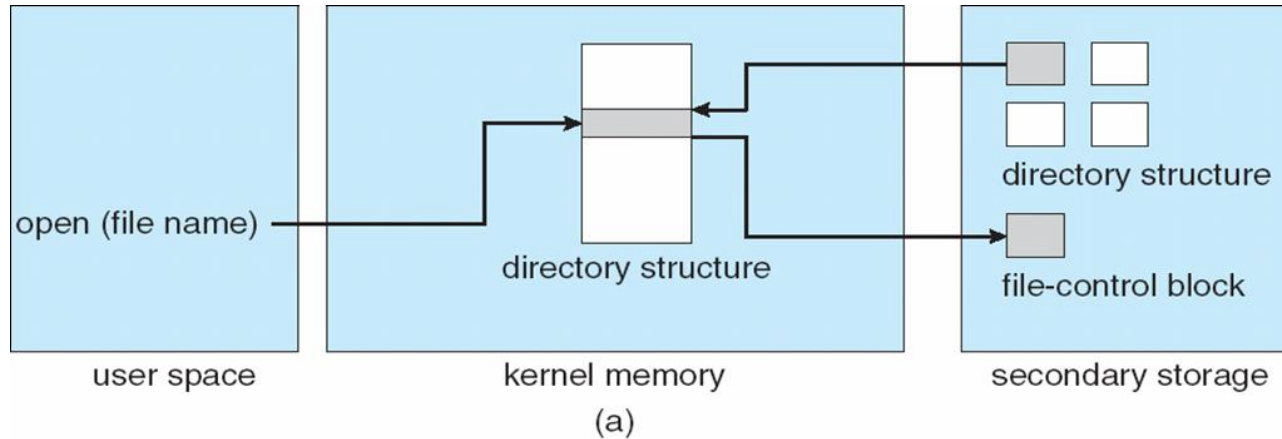| file permissions |
| --- |
| file dates (create, access, write) |
| file owner, group, ACL |
| file size |
| file data blocks or pointers to file data blocks |

# In-Memory File System Structures

☐ Mount table storing file system mounts, mount points, file system types

☐ The following figure illustrates the necessary file system structures provided by the operating systems

☐ Figure 12-3(a) refers to opening a file

☐ Figure 12-3(b) refers to reading a file

☐ Plus buffers hold data blocks from secondary storage

☐ Open returns a file handle for subsequent use

☐ Data from read eventually copied to specified user process memory address

# In-Memory File System Structures

# Partitions and Mounting

- Partition can be a volume containing a file system ("cooked") or **raw** – just a sequence of blocks with no file system

- Boot block can point to boot volume or boot loader set of blocks that contain enough code to know how to load the kernel from the file system

    - Or a boot management program for multi-os booting

- **Root partition** contains the OS, other partitions can hold other Oses, other file systems, or be raw

    - Mounted at boot time

    - Other partitions can mount automatically or manually

- At mount time, file system consistency checked

    - Is all metadata correct?

        ‣ If not, fix it, try again

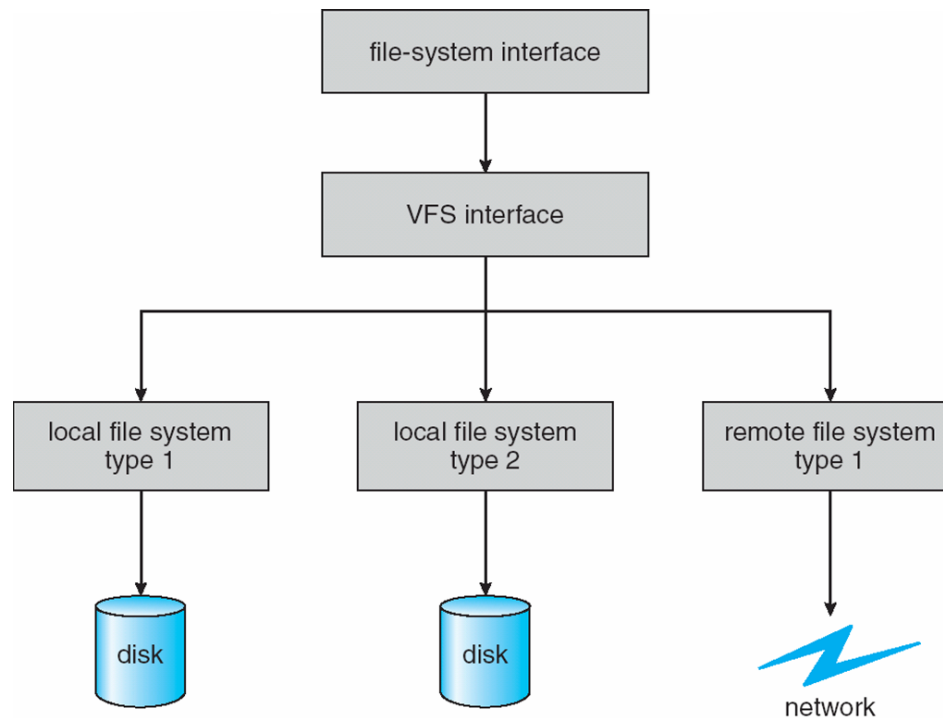        ‣ If yes, add to mount table, allow access

# Virtual File Systems

- **Virtual File Systems** (**VFS**) on Unix provide an object-oriented way of implementing file systems

- VFS allows the same system call interface (the API) to be used for different types of file systems

  - Separates file-system generic operations from implementation details

  - Implementation can be one of many file systems types, or network file system

    ▸ Implements **vnodes** which hold inodes or network file details

  - Then dispatches operation to appropriate file system implementation routines

# Virtual File Systems (Cont.)

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

# Virtual File System Implementation

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

# End