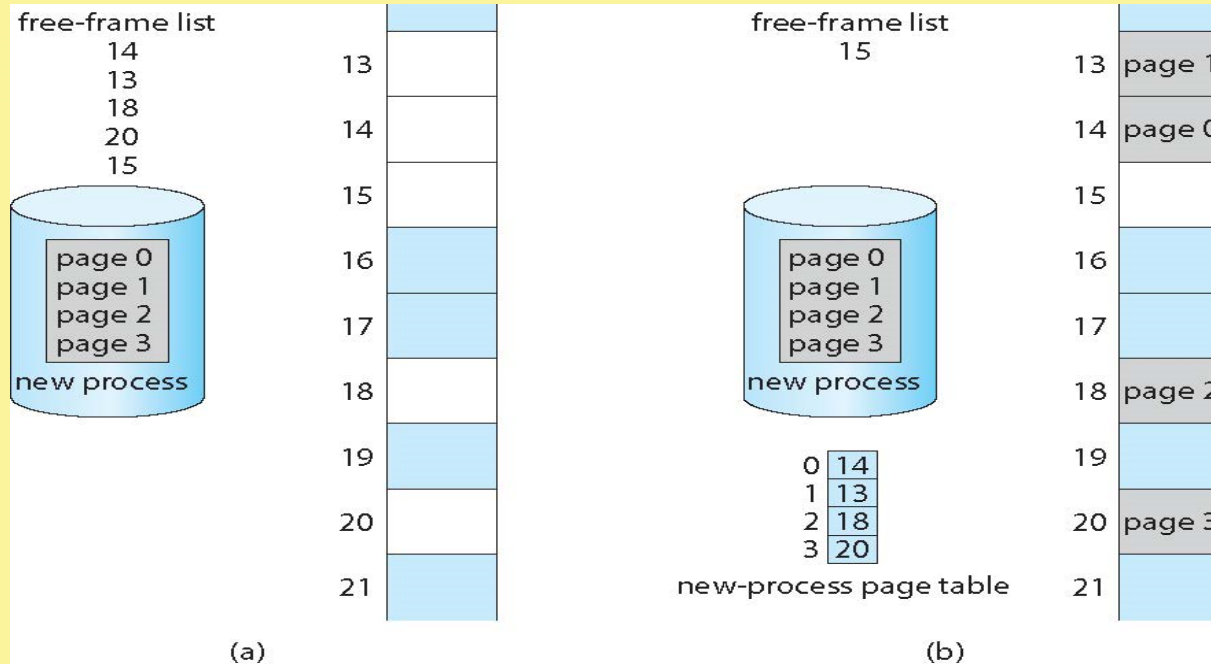# Allocating Frames to a New Process



Before allocation                    After allocation

# TLB -- Associative Memory

- If page table is kept in main memory 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 partially solved by the use of a special fast-lookup hardware cache called **associative memory** or **translation look-aside buffers** (**TLBs**)
- Associative memory – parallel search

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

- Address translation (p, d)
  - If p is in associative register, get frame # out
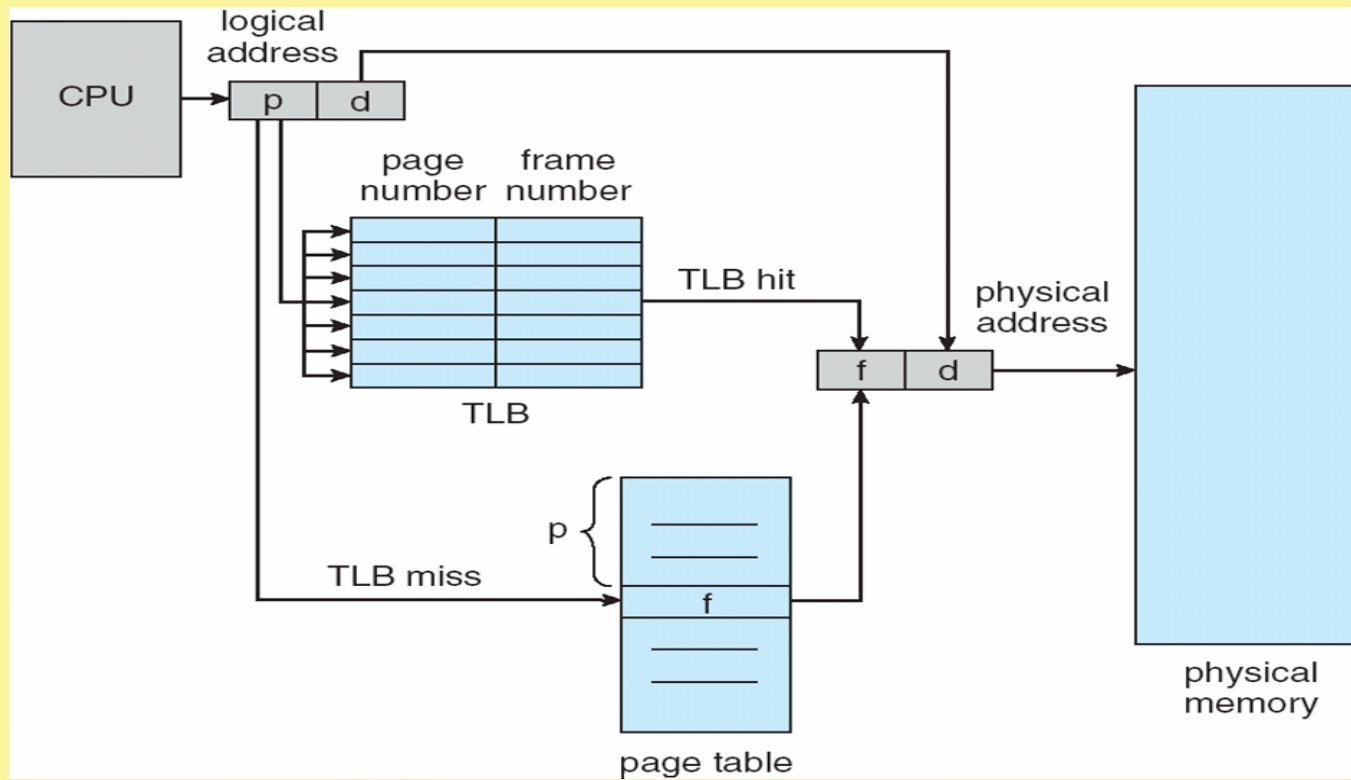  - Otherwise get frame # from page table in memory

# TLB issues

- TLB is typically small (64 to 1,024 entries)
- On a TLB miss, the value of the (missed page-table and frame-number), is loaded into the TLB for faster access next time that address is used.
  - What if there is no free TLB entry? Replacement policies must be considered
  - Some entries can be **wired down** for permanent fast access
- 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 TLB at every context switch

# 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
- **Effective Access Time** (**EAT**)

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

$$= 2 + \varepsilon - \alpha$$

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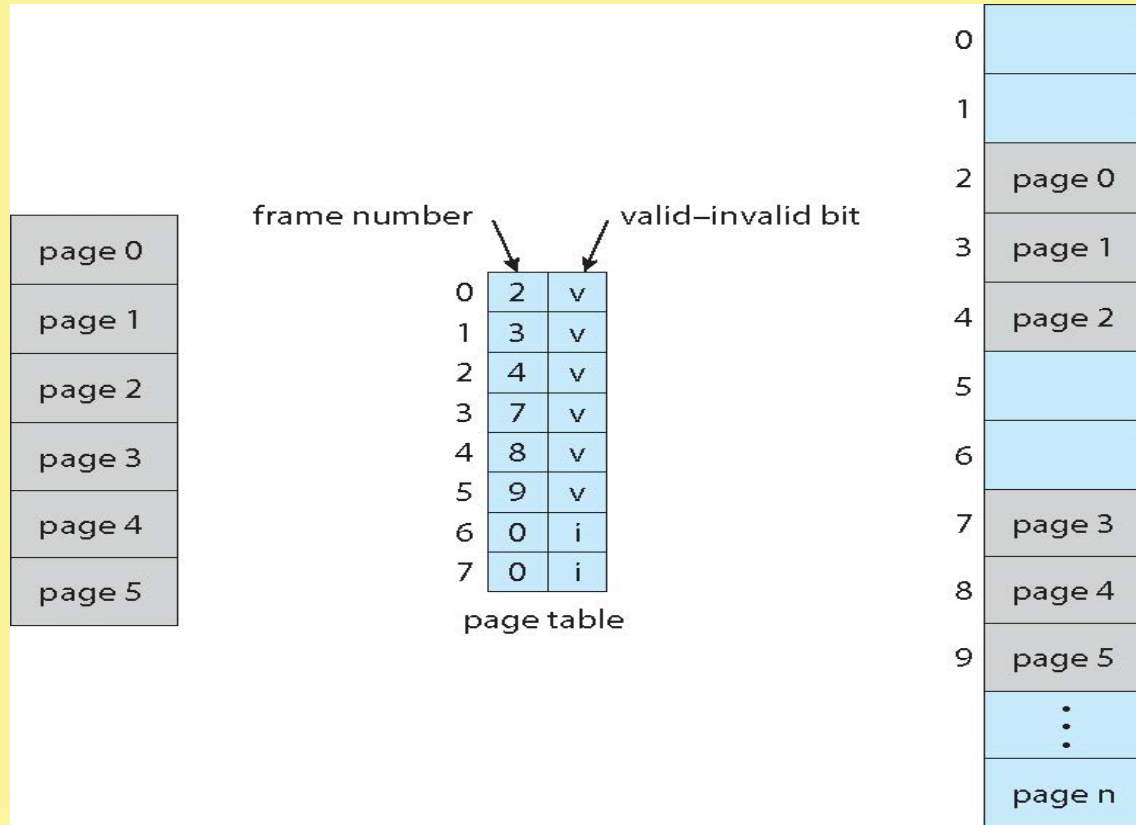    - EAT = 0.99 x 100 + 0.01 x 200 = 101ns

# Memory Protection

- Memory protection implemented by associating protection bits with each frame to indicate if "read-only " or "read-write" access is allowed
  - Can also add more bits to indicate "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 is 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



page 0
page 1
page 2
page 3
page 4
page 5

frame number          valid–invalid bit

| | | |
|---|---|---|
| 0 | 2 | v |
| 1 | 3 | v |
| 2 | 4 | v |
| 3 | 7 | v |
| 4 | 8 | v |
| 5 | 9 | v |
| 6 | 0 | i |
| 7 | 0 | i |

page table

0
1
2  page 0
3  page 1
4  page 2
5
6
7  page 3
8  page 4
9  page 5
⋮
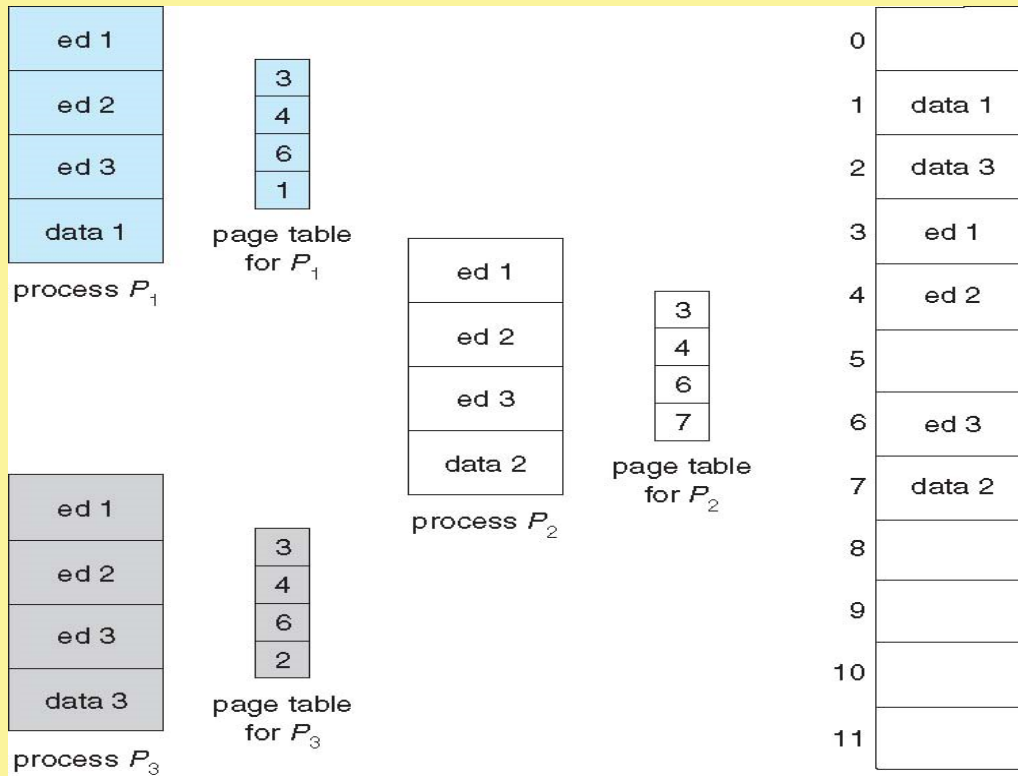page n

# 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 inte-rprocess 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
  - Page size of 1 KB ($2^{10}$)
  - Page table would have 4 million entries ($2^{32} / 2^{10}$)
  - If each entry is 4 bytes -> Page table is of size 16 MB
    - That amount of memory used to cost a lot.
    - Do not want to allocate that contiguously in main memory
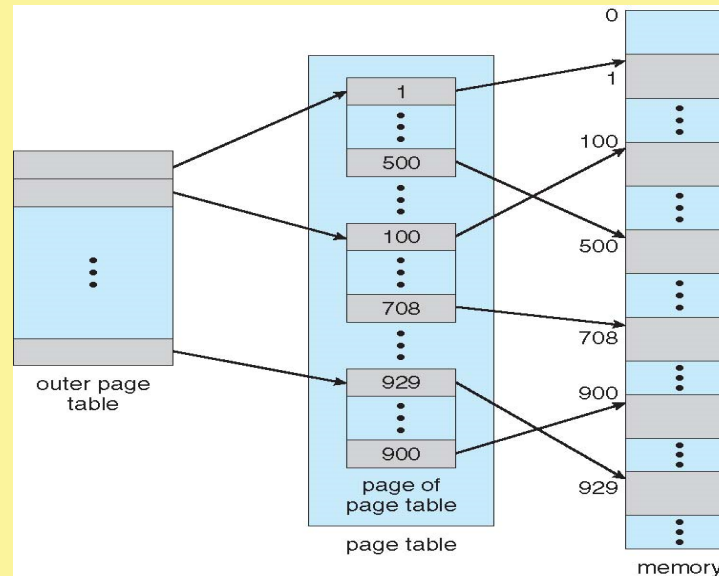- What about a 64-bit logical address space?

Page Table for Large address space

- 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

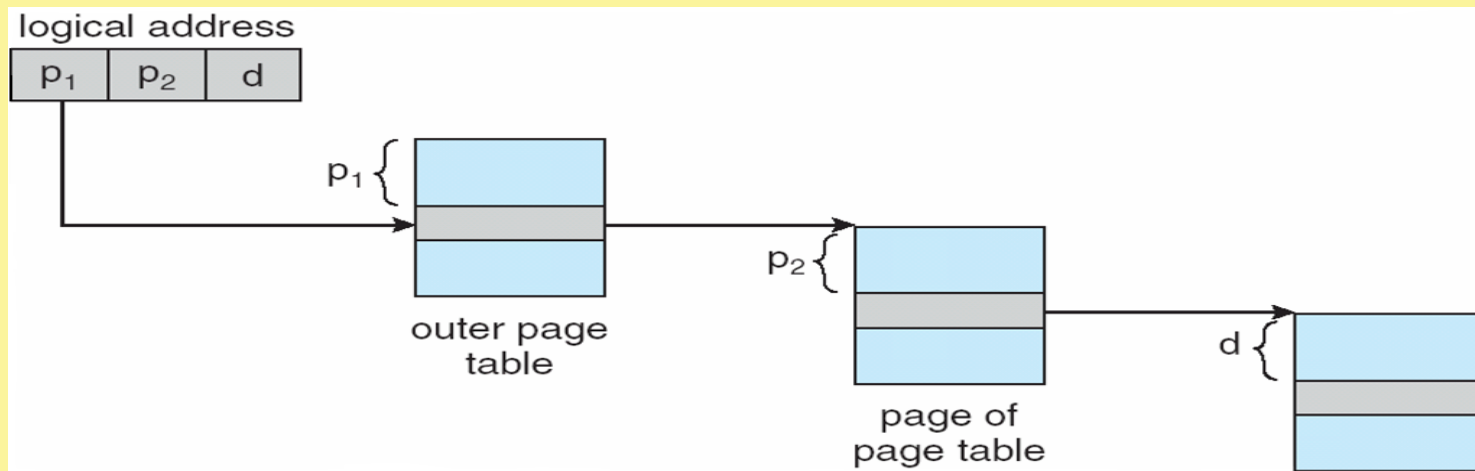# Two-Level Paging Example

- 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:

| $p_1$ | $p_2$ | $d$ |
|-------|-------|-----|
| 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

# 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 | offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

  - Outer page table has $2^{42}$ entries or $2^{44}$ bytes

# 64-bit Logical Address Space (Cont.)

- One solution is to divide the outer page table. Various ways of doing so. Example – three-level page table

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

- Even with 2nd outer page table, the outer-outer table is still $2^{34}$ bytes in size.
- And possibly 4 memory access to get to one physical memory location.
- The next step would be four-level. But ….

# 64-bit Logical Address Space (Cont.)

Several schemes for dealing with very large logical address space

- Hashed Page Table.

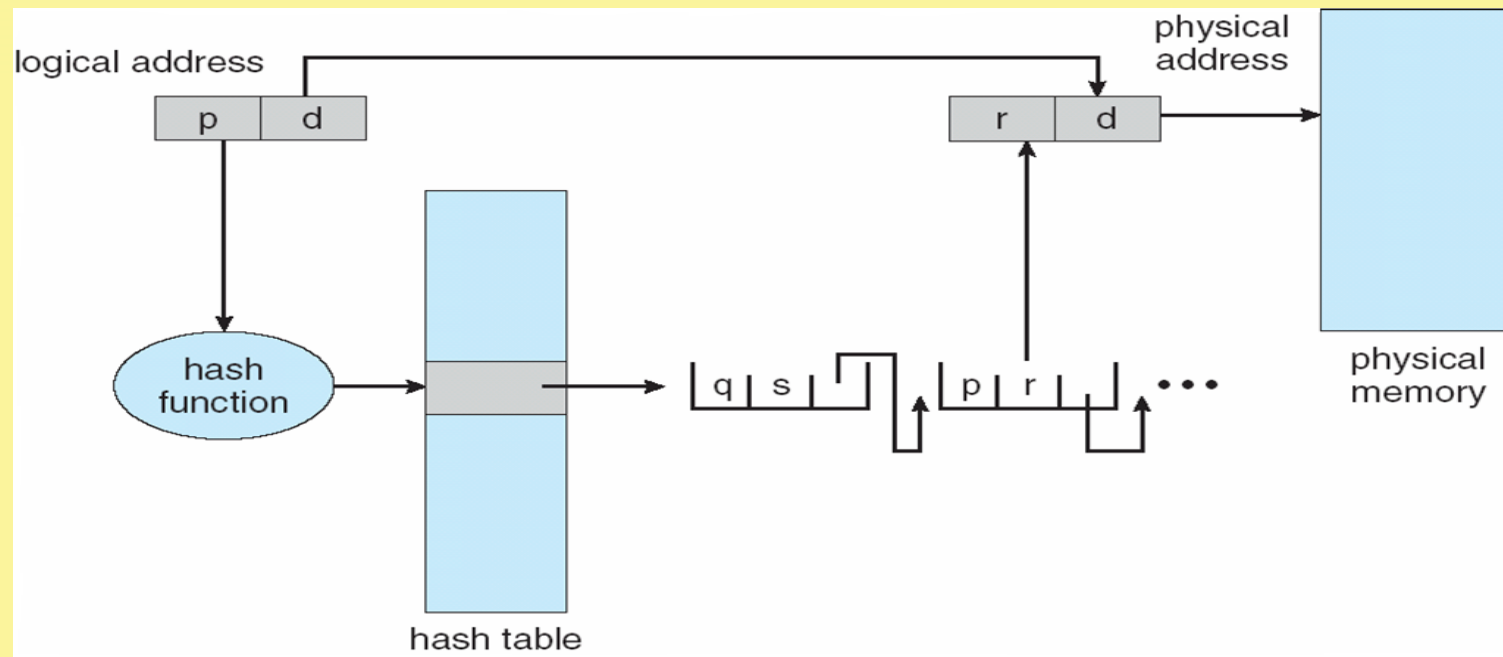- Clustered Page Tables

- Inverted Page Table

# Hashed Page Table

- 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
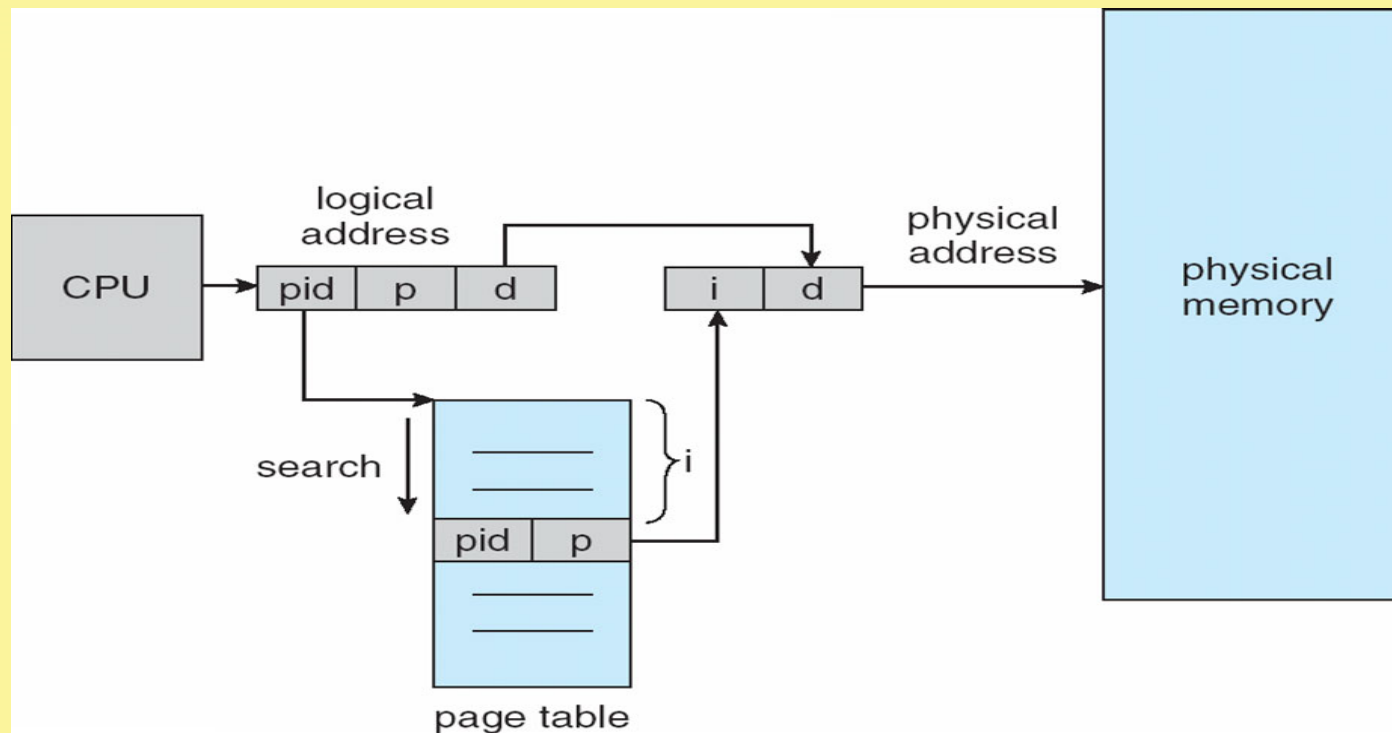
# Hashed Page Table

# Inverted Page Table

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

- Use **inverted page-table**, which has one entry for each real page of memory

- An entry the inverted-page table consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page.

- What is maximum size of the inverted page-table?

# Inverted Page Table Architecture

# Inverted Page Table (Cont.)

- Decreases memory needed to store each individual page table, but increases time needed to search the inverted page 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

# Oracle SPARC Solaris

- Consider a 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.
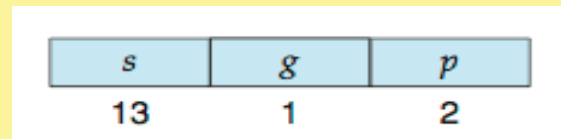
Example: The Intel IA-32 Architecture

- 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**))

- 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
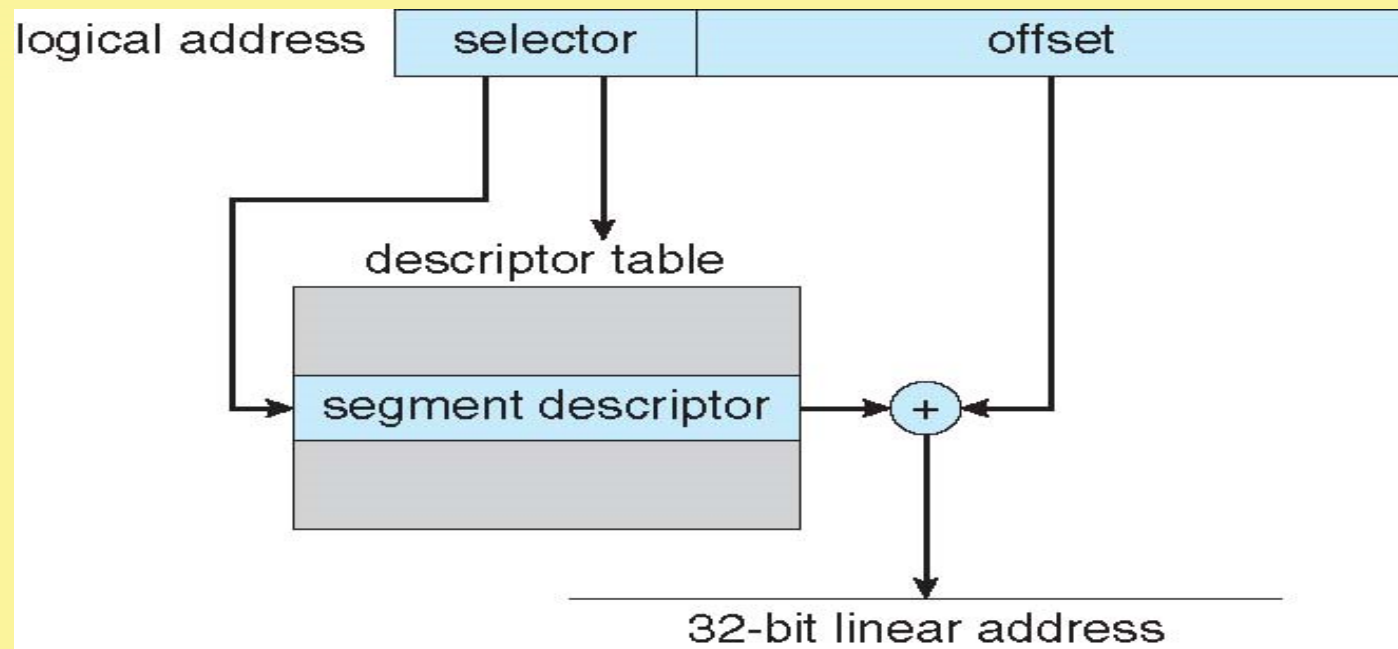    - Pages sizes can be 4 KB or 4 MB
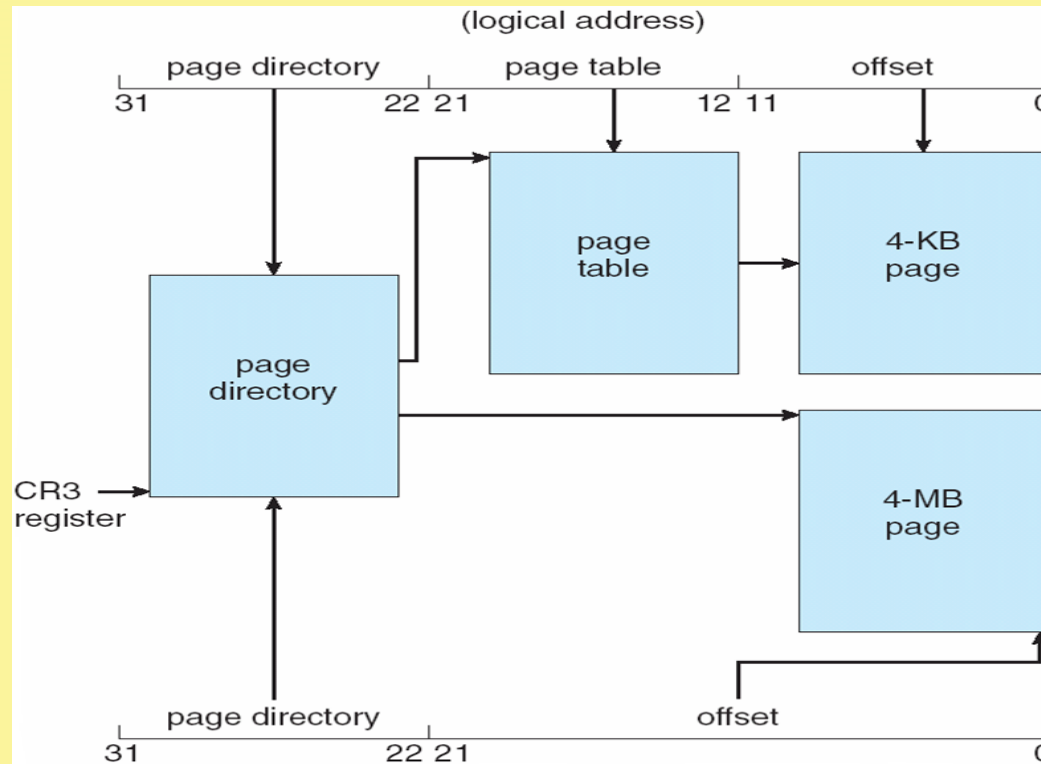
# Logical to Physical Address Translation in IA-32



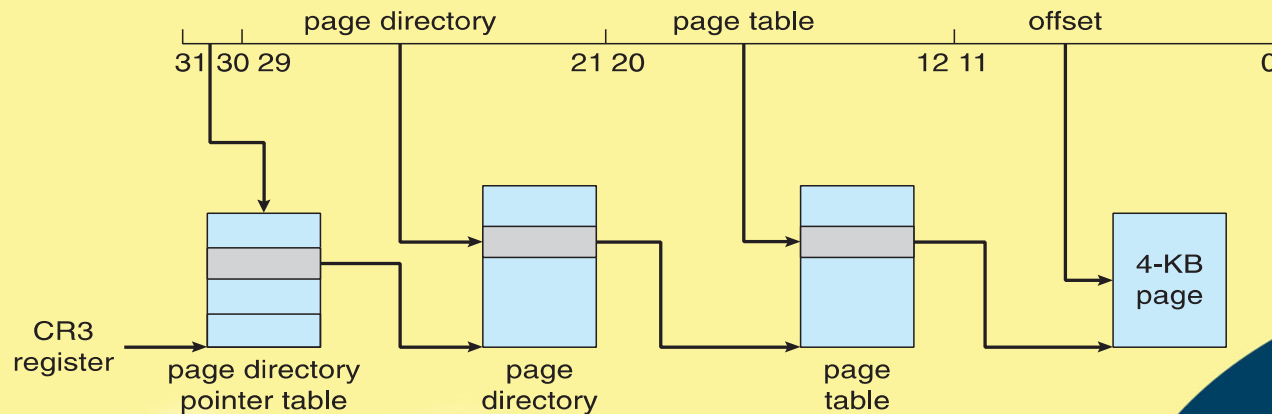| page number | | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 10 | 10 | 12 |

# Intel IA-32 Segmentation
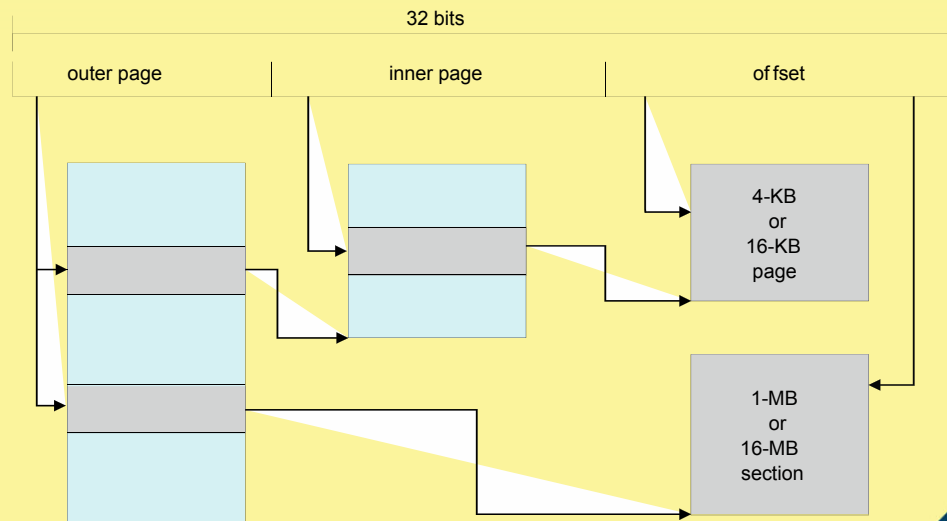
# Intel IA-32 Paging Architecture

# 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

# 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

*Conclusion*

**Conclusion:**

- **Memory management is an important OS activity**

- **Various strategies have been discussed**

- **Paging appears to be a good approach**

- **Segmentation provides a logical division of program**

- **Advanced processors support a mix of many of the**

  **strategies**

**NPTEL ONLINE CERTIFICATION COURSES**

Thank you

**Operating System Fundamentals**
**Santanu Chattopadhyay**
**Electronics and Electrical Communication Engg.**

# Virtual Memory

# Concepts Covered:

- Background
- Demand Paging
- Copy-on-Write
- Page Replacement
- Allocation of Frames
- Thrashing
- Memory-Mapped Files
- Other Considerations

# Objectives

- To describe the benefits of a virtual memory system

- To explain the concepts of demand paging, page-replacement algorithms, and allocation of page frames

- To discuss the principle of the working-set model

- To examine the relationship between shared memory and memory-mapped files

- To explore how kernel memory is managed

# Background

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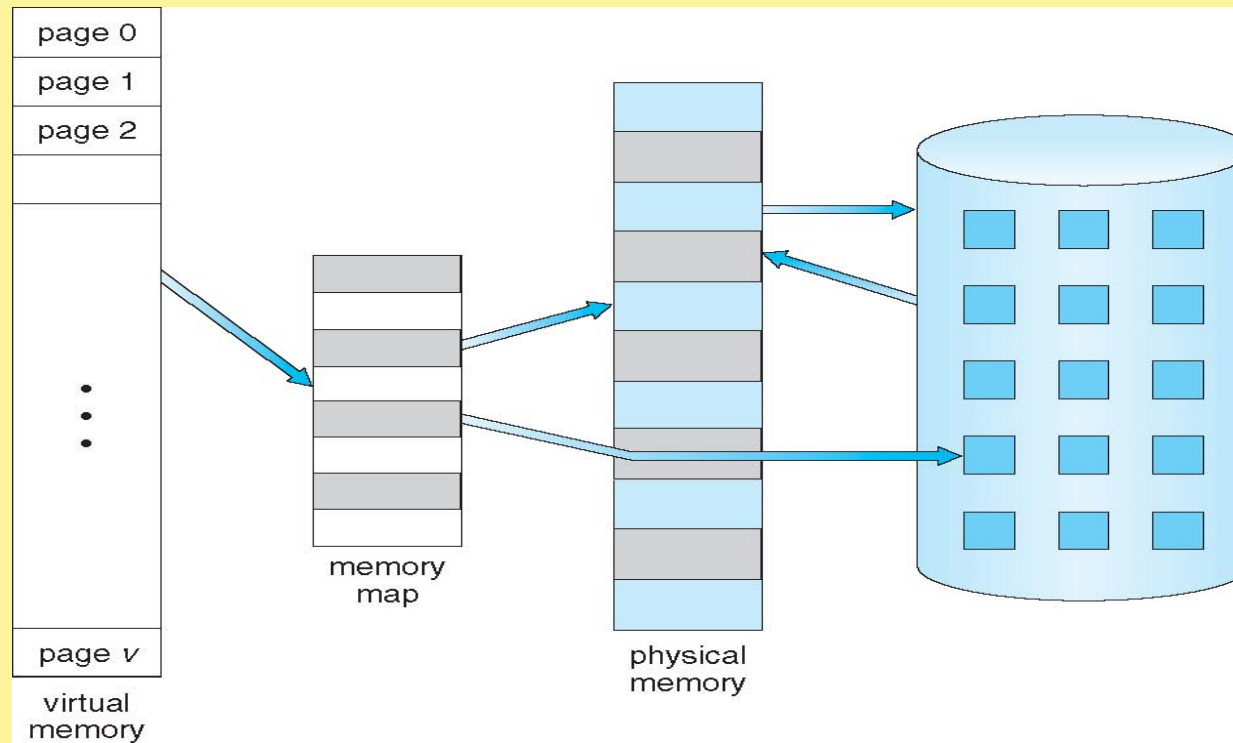
# Virtual Memory

**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:
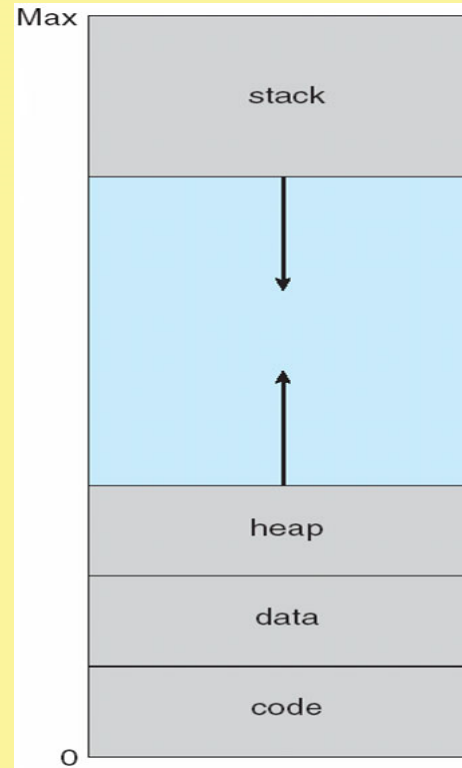  - Demand paging
  - Demand segmentation

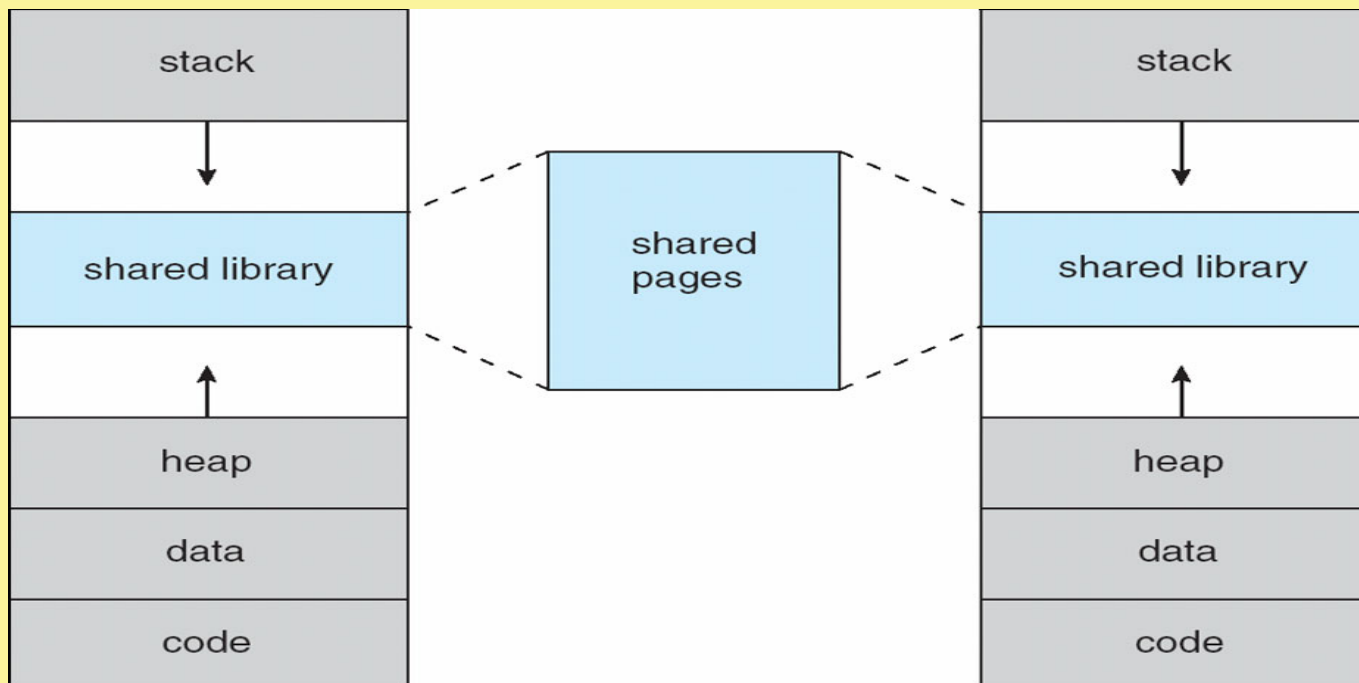# Virtual Memory That is Larger Than Physical Memory

# Virtual-address Space

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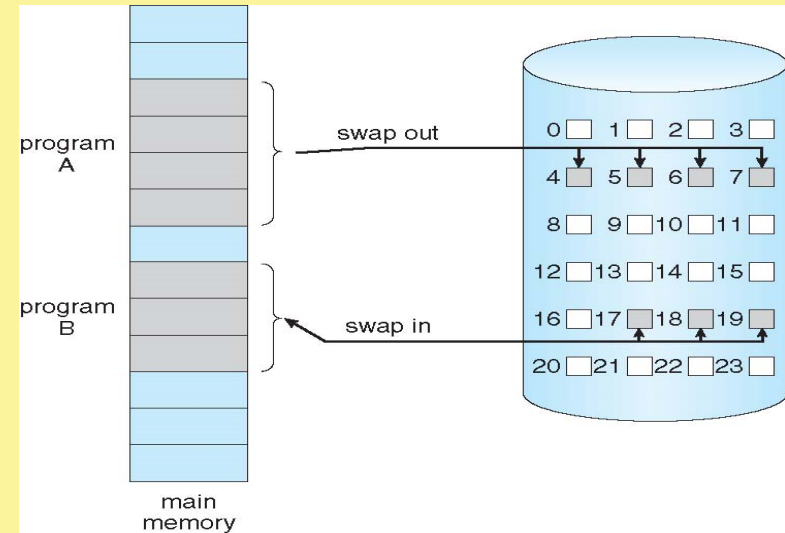- Pages can be shared during `fork()`, speeding process creation

# Shared Library Using Virtual Memory

# Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed
  - Less I/O needed, no unnecessary I/O
  - Less memory needed
  - Faster response
  - More users
- Similar to paging system with swapping (diagram on right)
- Page is needed ⇒ reference to it
  - invalid reference ⇒ abort
  - not-in-memory ⇒ bring to memory
- **Lazy swapper** – never swaps a page into memory unless page will be needed
  - Swapper that deals with pages is a **pager**

# Basic Concepts

- When a process is to be swapped in, the pager guesses which pages will be used before the process is swapped out again

- Instead of swapping in a whole process, the pager brings in only those "guessed" pages into memory

- Need new MMU functionality to implement demand paging. Need to distinguish between the pages that are in memory and the pages that are on the disk. Use a variation of the valid-invalid scheme used for protection (see next slide)

- If pages needed are already **memory resident**
  - No difference from non demand-paging

- If page needed and not memory resident, need to detect and load the page into memory from storage
  - Without changing program behavior
  - Without programmer needing to change code

# Valid-Invalid Bit

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

- Initially valid–invalid bit is set to $i$ on all entries

- Example of a page table snapshot shown here

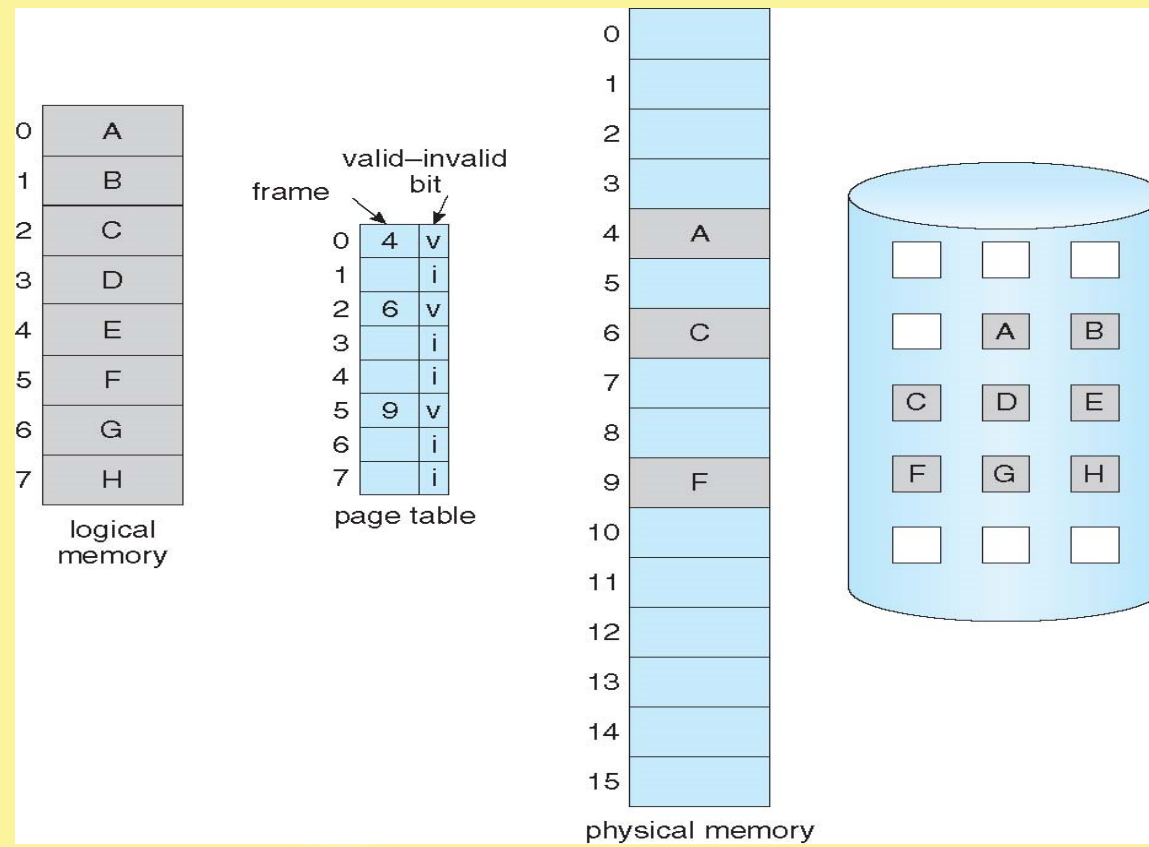| Frame # | valid-invalid bit |
|---------|-------------------|
|         |                   |
|         | v                 |
|         | v                 |
|         | v                 |
|         | i                 |
| . . .   |                   |
|         | i                 |
|         | i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

page table

- During MMU address translation, if valid–invalid bit in page table entry is $i \Rightarrow$ page fault

44

## 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**

Actions when a page fault occurs

1.Operating system looks at another table to decide:
- Invalid reference $\Rightarrow$ abort
- Just not in memory. Go to step (2).

2.Find free frame

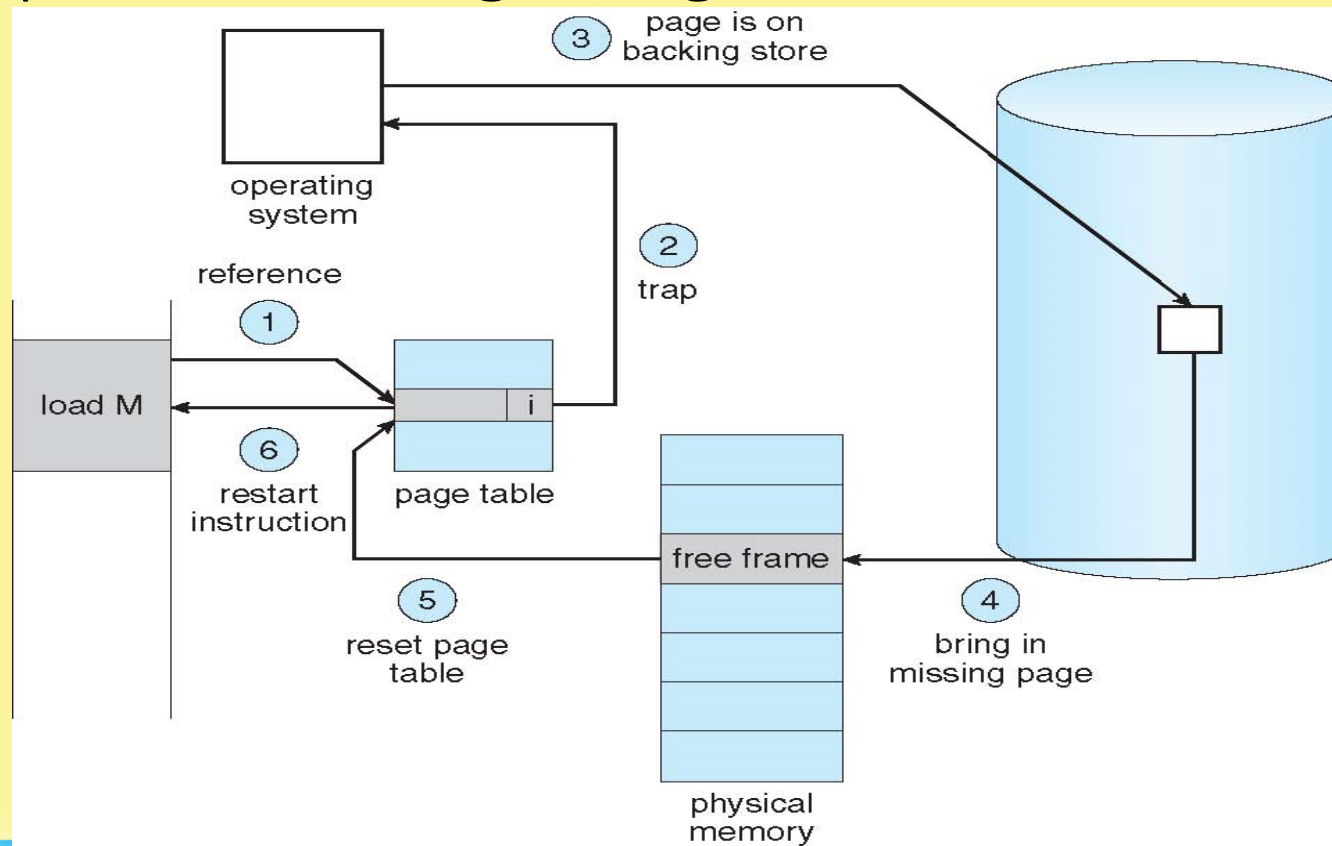3.Swap page into frame via scheduled disk operation

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

5.Restart the instruction that caused the page fault
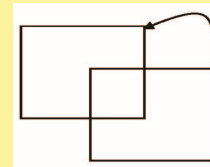
# Steps in Handling a Page Fault

# Aspects of Demand Paging

- **Pure demand paging** – start process with *no* pages in memory
  - OS sets instruction-pointer to the first instruction of the process, non-memory-resident -> page fault
  - And for every other process pages on first access
- Actually, a given instruction could access multiple pages → multiple page faults
  - Consider fetch and decode of instruction which adds 2 numbers from memory and stores result back to memory
  - The two numbers may reside in two different pages
- 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
  - Page fault during the move
  - Restart the whole operation?

- What if source and destination overlap?



- Several solution.  Simplest one is to figure out all the pages that are needed to execute the instruction to completion and ensure that these pages are in memory before the instruction starts executing.

# Stages in Demand Paging

Stages in Demand Paging (worse case)

1. Trap to the operating system

2. Save the user registers and process state

3. Determine that the interrupt was a page fault

4. Check that the page reference was legal and determine the location of the page on the disk

5. Issue a read I/O 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

# Stages in Demand Paging (Cont.)

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