# CS370 Operating Systems

**Colorado State University**

**Yashwant K Malaiya**

**Fall 22 L18**

**Main Memory**

OS is a *systems* class, where hardware and software come together.

**Slides based on**
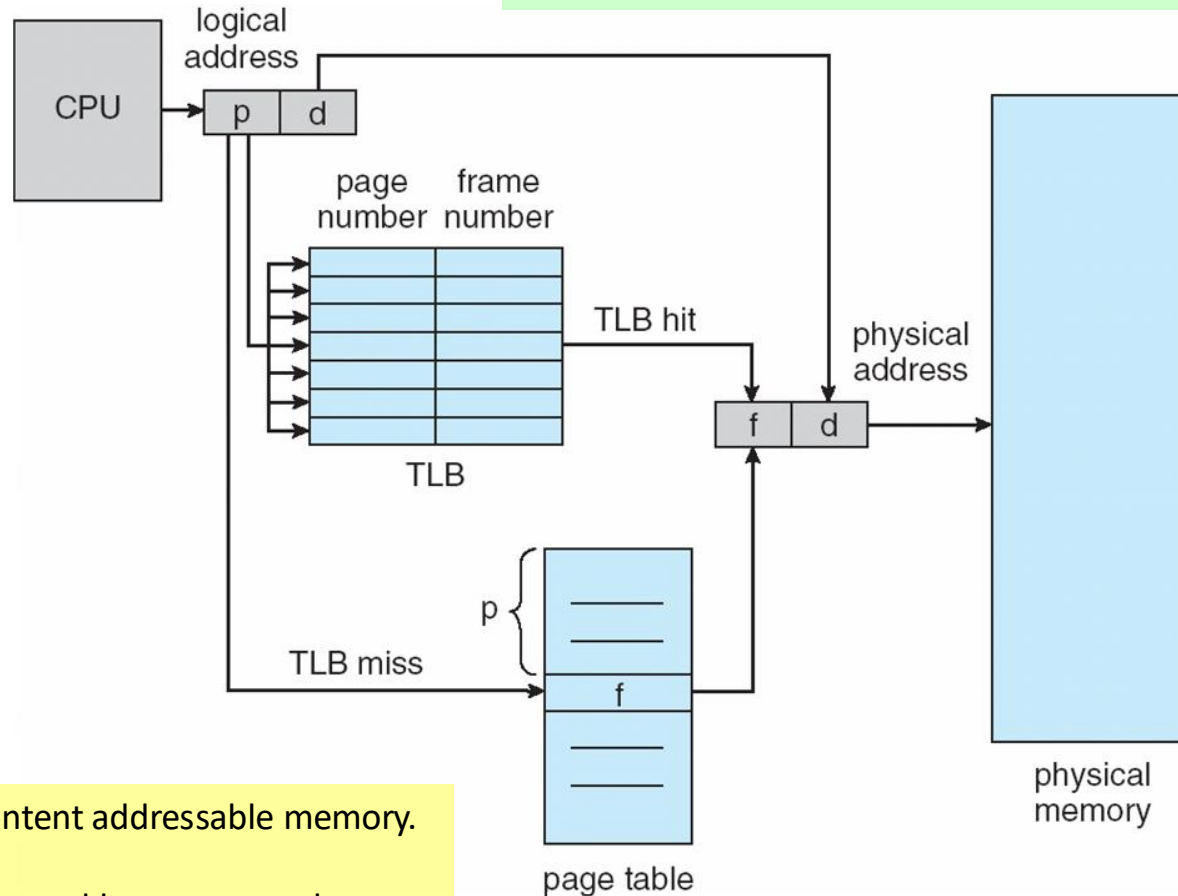- Text by Silberschatz, Galvin, Gagne
- Various sources

# FAQ

- Why use pages? So that memory does not have be allocated contiguously.

- Where is the page table? Memory, with a part cached in TLB

- How to find the page table in memory? Page table base register

- Is there is specific formula for calculating the physical address from the logical address? Page number to frame number lookup

- Each process has its own page table? Can there be a conflict in sharing physical memory? No, unless..

- Where is the TLB ? On the same chip as CPU.

- Why use associative memory for TLBs? Fast content-based search to find frame number

# Paging Hardware With TLB

Page number  p  to frame number   f



TLB: uses content addressable memory.

TLB Miss: page table access may be done using hardware or software

Colorado State University

# Effective Access Time

**General approach:**   expected access time
Effective access time
$$= Pr\{\text{access type A}\}. \text{Access-time}_A +$$
$$Pr\{\text{access type B}\}. \text{Access-time}_B$$

**Ex: effective access time with TLB/page table:**
- Associative Lookup = $\varepsilon$ time units
- Hit ratio = $\alpha$
- **Effective Access Time** (**EAT**): probability weighted
  $$EAT = (100 + \varepsilon)\, \alpha + (200+\varepsilon)(1 - \alpha)$$
- Ex:
  Consider $\alpha = 80\%$, $\varepsilon$ = negligible for TLB search, 100ns for memory access
  – EAT = 100x0.80 + 200x0.20 = 120ns

**Colorado State University**

Colorado State University

# 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
- Any violations result in a trap to the kernel

Colorado State University

"invalid" : page is not in the process's address space.

# Shared Pages among Processes

- **Shared code**
  - One copy of read-only (**reentrant** non-self modifying) 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

**Colorado State University**

# Shared Pages Example



ed1, ed2, ed3
(frames 3, 4, 6) shared

Colorado State University

Optimal Page Size Computation:

  page table size vs internal  fragmentation tradeoff

- Average process size = *s*

- Page size = *p*

- Size of each entry in page table = *e*
  - Pages per process = *s/p*
  - *se/p:* Total page table space for average process
  - Total Overhead = Page table overhead + Internal fragmentation loss
    
    = *se/p* + *p/2*

Colorado State University

- Total Overhead = $se/p + p/2$
- Optimal: Obtain derivative of overhead with respect to $p$, equate to 0

  $-se/p2 + 1/2 = 0$

- i.e.    $p^2 = 2se$    or $p = (2se)^{0.5}$

**Assume   $s$ = 128KB and $e=8$ bytes per entry**

- Optimal page size = 1448 bytes
  - In practice we will never use 1448 bytes
  - Instead, either 1K or 2K would be used
    - **Why?** Pages sizes are in powers of 2 i.e. $2^X$
    - Deriving offsets and page numbers is also easier

**Colorado State University**

# Page Table Size

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

- Consider a 32-bit logical address space as on recent processors 64-bit on 64-bit processors
  - Assume page size of 4 KB ($2^{12}$) entries
  - 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
    - Don't want to allocate that **contiguously** in main memory

| $2^{10}$ | **1024  or 1 kibibyte** |
|---|---|
| $2^{20}$ | 1M   mebibyte |
| $2^{30}$ | 1G    gigibyte |
| $2^{40}$ | 1T    tebibyte |

**Colorado State University**
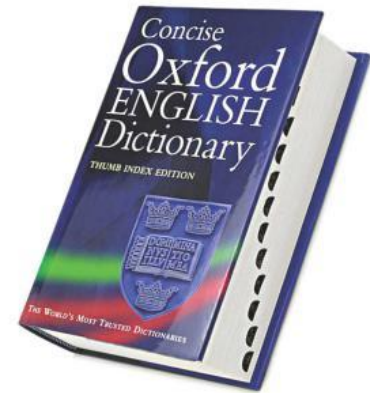
# Issues with large page tables

- Cannot allocate page table **contiguously** in memory
- Solution:
  - Divide the page table into smaller pieces
  - **Page the page-table**
    - Hierarchical Paging

**Colorado State University**

# Hierarchical Page Tables

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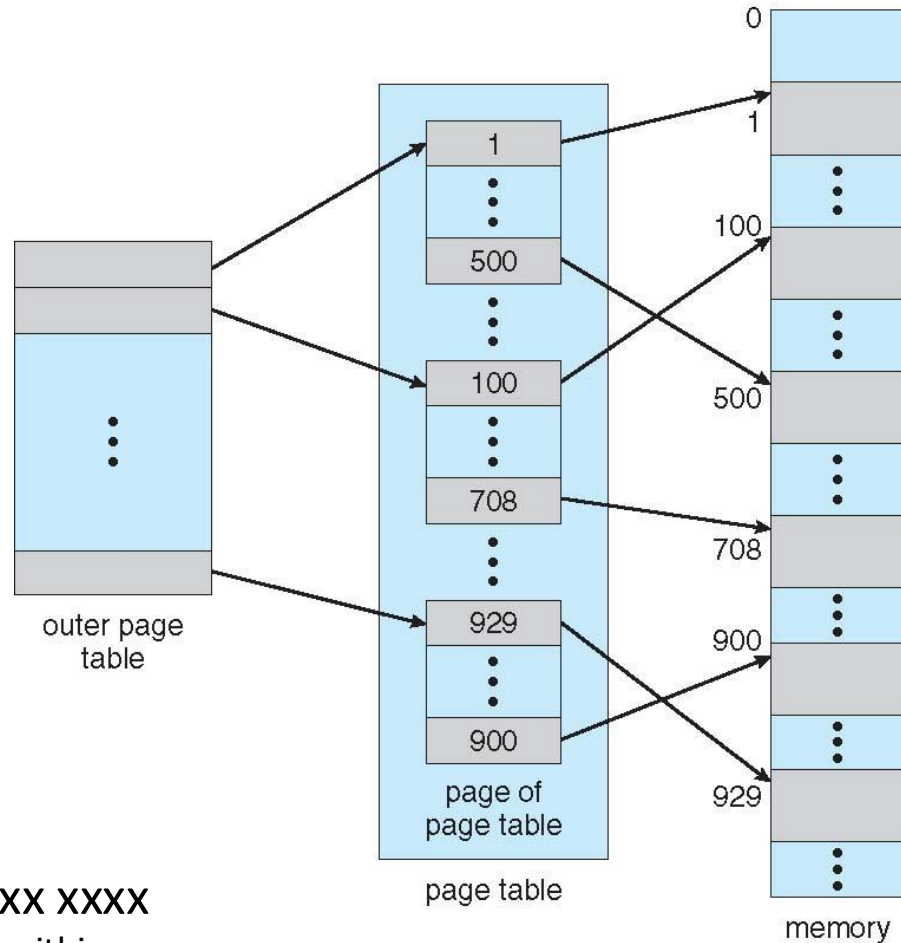- We then page the page table

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

P1: indexes the outer page table
P2: page table: maps to frame

Colorado State University

# Two-Level Page-Table Scheme

page number      page offset

| $p_1$ | $p_2$ | $d$ |
|-------|-------|-----|
| 12 | 10 | 10 |

outer page table

page of page table

page table

memory

XXXX XXXX XXXX XXXX XXXX XX XX XXXX XXXX

Outer Page table      page table      offset within page

**Colorado State University**

# 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**

**Colorado State University**
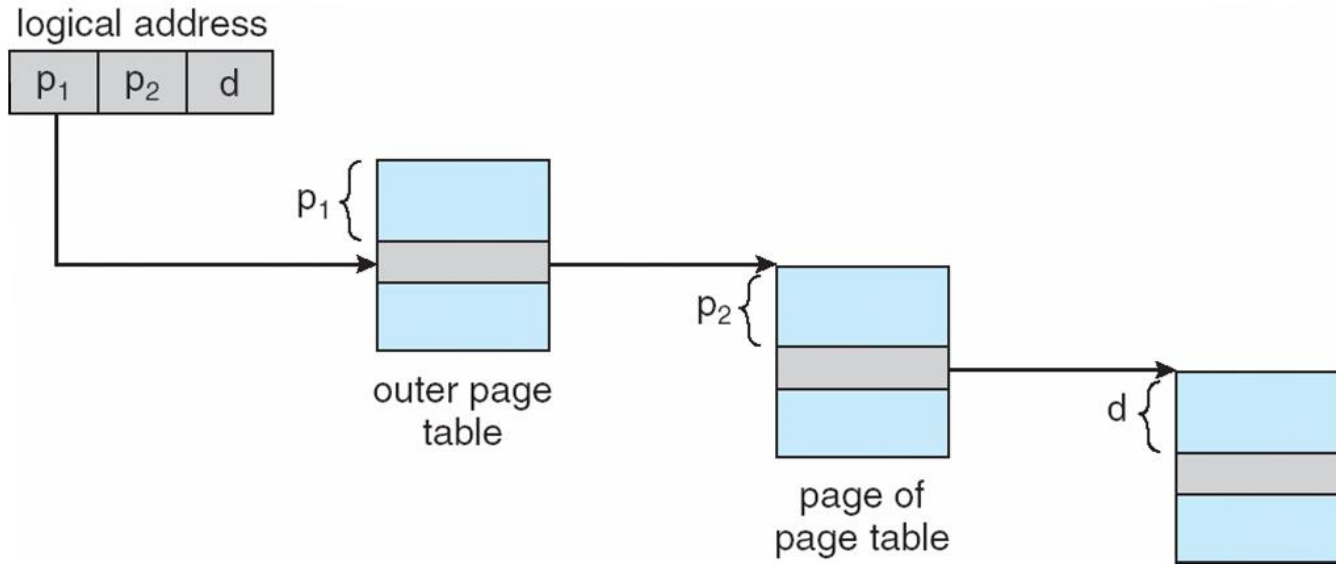
# Two-Level Paging Example

- A logical address is as follows:

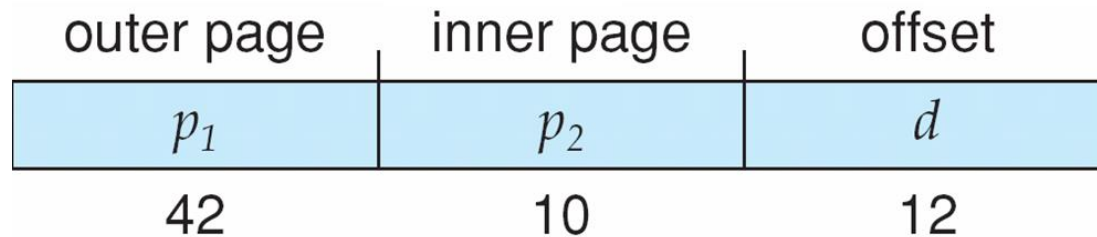| page number | | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 12 | 10 | 10 |

- One Outer page table: size $2^{12}$
  entry: page of the page table

- Often only some of all possible $2^{12}$ Page tables needed (each of size $2^{10}$)
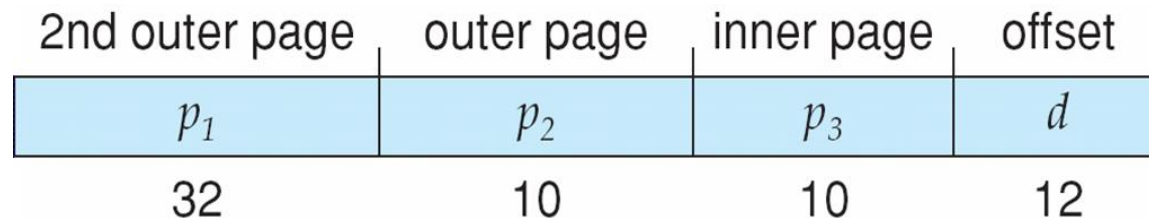
**Colorado State University**

# Hierarchical Paging



If there is a hit in the TLB (say 95% of the time), then average access time will be close to slightly more than one memory access time.

**Colorado State University**

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

- Problem: Outer page table has $2^{42}$ entries!
- Approach: Divide the outer page table into 2 levels
  - 4 memory accesses!

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

**Colorado State University**

# Three-level Paging Scheme

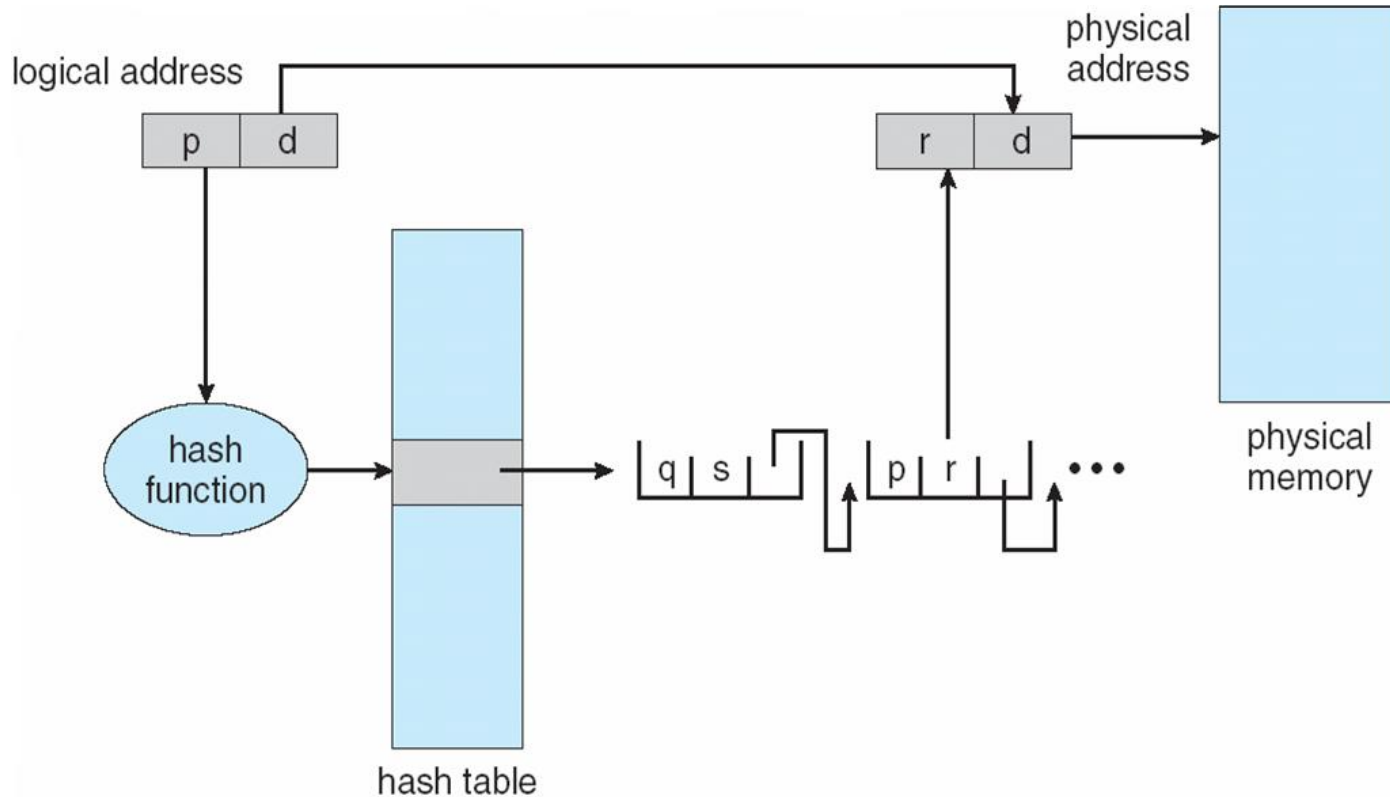| outer page | inner page | offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

- Outer page table has $2^{42}$ entries!
- Divide the outer page table into 2 levels
  - 4 memory accesses!

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

**Colorado State University**

# Hashed Page Tables

- Useful when 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)
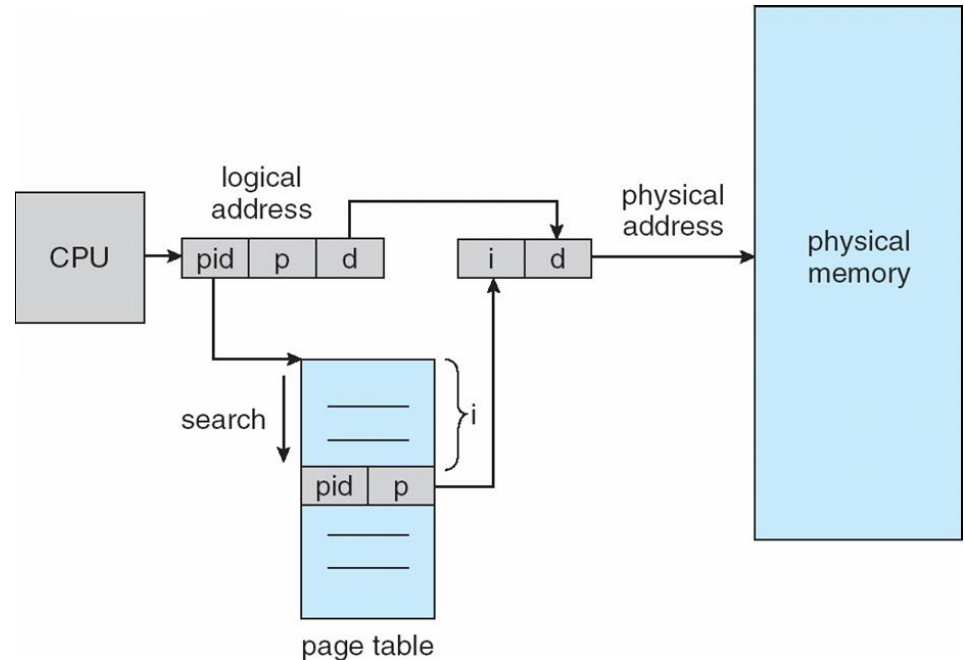
**Colorado State University**

# Hashed 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

Colorado State University

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

Search for pid, p, offset i is the physical frame address
Note: multiple processes in memory
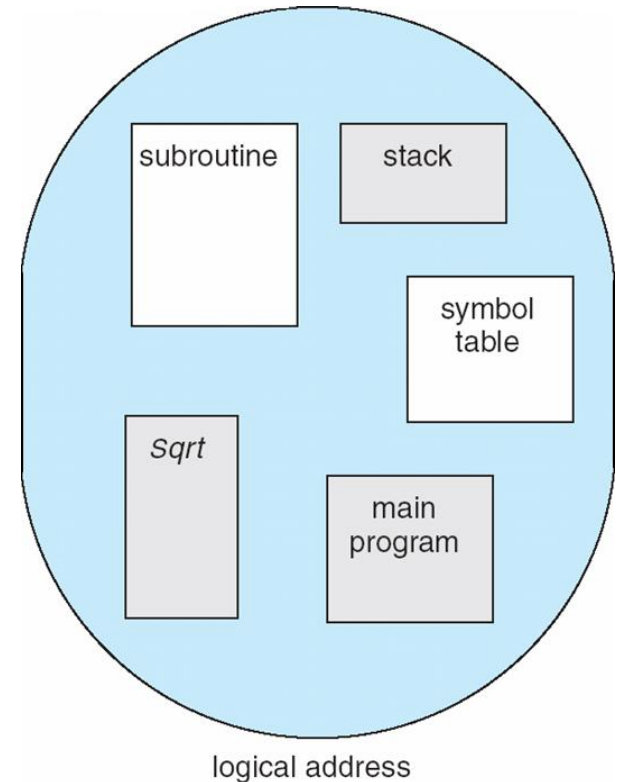
# Inverted Page Table

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

- But how to implement shared memory?
  - One mapping of a virtual address to the shared physical address. Not possible.

Used in IA-64 ..

Colorado State University

# Segmentation Approach

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, arrays, symbol table

- Segment table
  - Segment-table base register (STBR)
  - Segment-table length register (STLR)
- segments vary in length, can very dynamically
- Segments may be paged
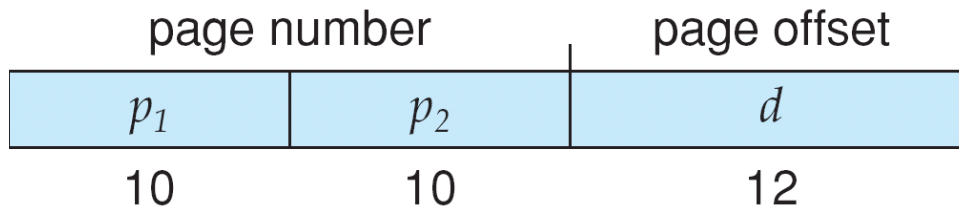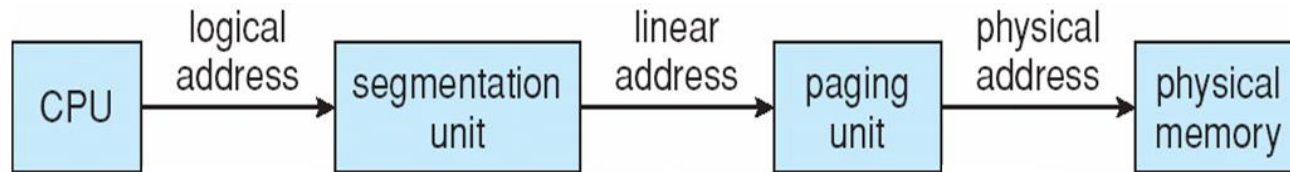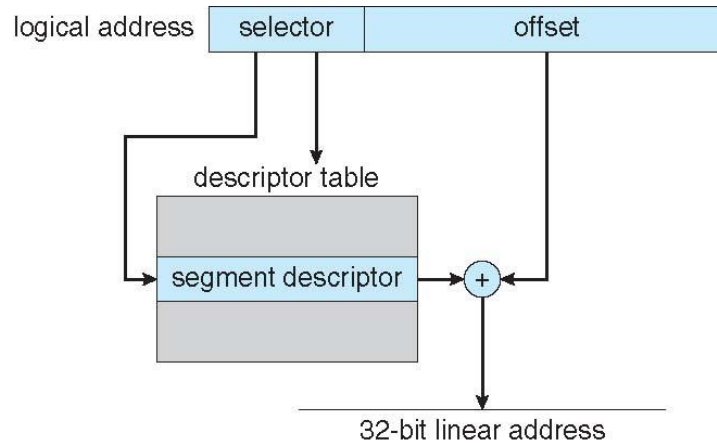- Used for x86-32 bit
- Origin of term "segmentation fault"

subroutine

stack

symbol table

Sqrt

main program

logical address

**Colorado State University**

- Intel IA-32 (x386-Pentium)
- x86-64 (AMD,  Intel)
- ARM (Acorn > ARM Ltd > Softbank > ~~Nvidea~~)

Market: Upward compatibility.

Question: Why don't all the designers all use one  single approach?

**Colorado State University**

logical address | selector | offset

descriptor table

segment descriptor + 

32-bit linear address

CPU → logical address → segmentation unit → linear address → paging unit → physical address → physical memory

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

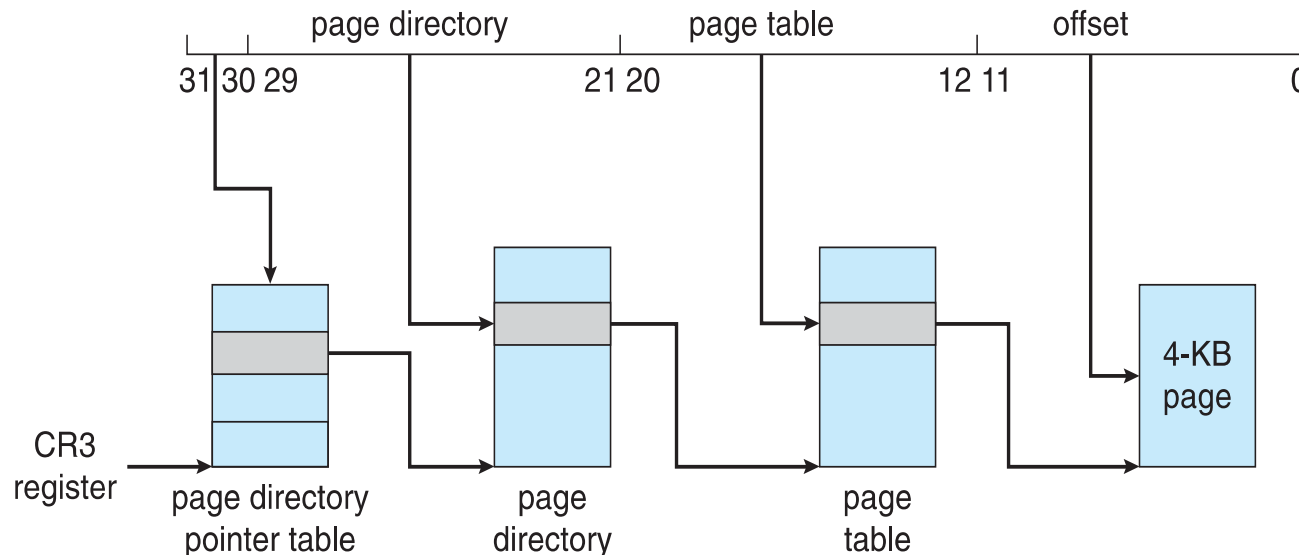**Colorado State University**

27

# Intel IA-32 Paging Architecture
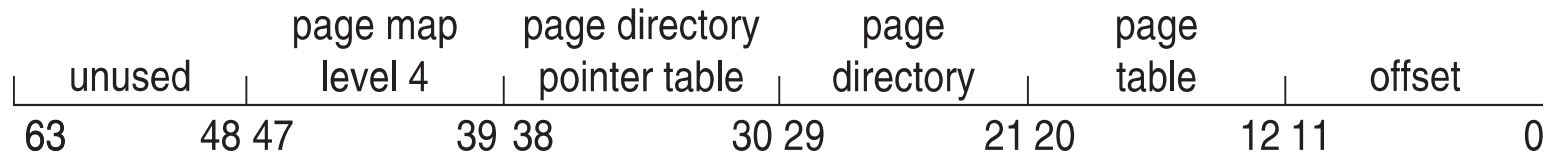


Support for two page sizes

# Intel IA-32 Page Address Extensions

n   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
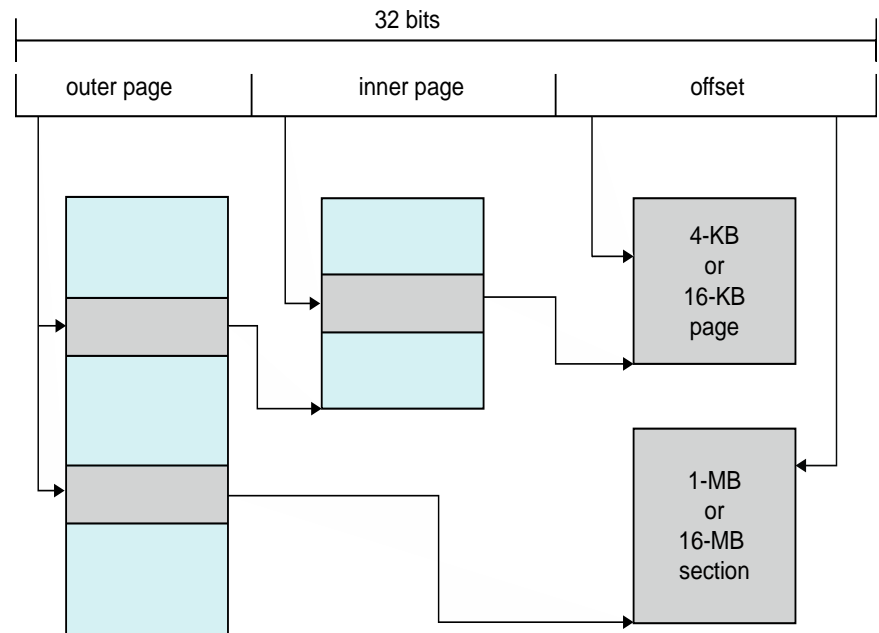- Net effect is increasing address space by increasing frame address bits.

# Intel x86-64

n   Intel x86 architecture based on AMD 64 bit architecture

n   64 bits is ginormous (> 16 exabytes)

n   In practice only implement 48 bit addressing or perhaps 52 or 57

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

  l   Four levels of paging hierarchy

n   Can also use PageAddressExtensions so virtual addresses are 48 bits and physical addresses are 52 (now 57)  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 |

Exabyte: $1024^6$ bytes

# Example: ARM Architecture

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

n    Modern, energy efficient, 32-bit CPU  now 64 bit also

n    4 KB and 16 KB pages

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

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

n    Two levels of TLBs

     l    Outer level has two micro TLBs (one data, one instruction)

     l    Inner is single main TLB

     l    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

**Colorado State University**

# CS370 Operating Systems

## Colorado State University
## Yashwant K Malaiya
## Fall 2022

## Virtual Memory

**Slides based on**
- Text by Silberschatz, Galvin, Gagne
- Various sources
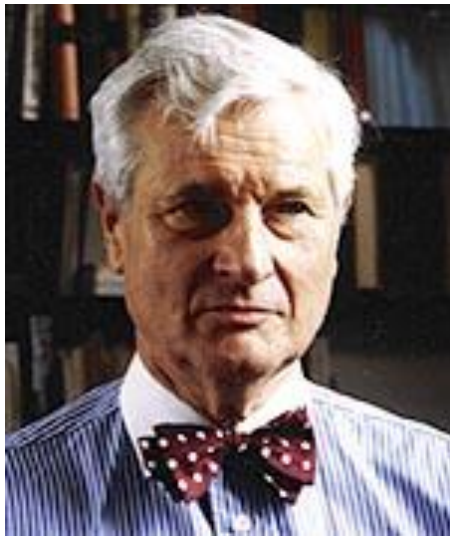
# Virtual Memory: Objectives



- A virtual memory system
- Demand paging, page-replacement algorithms, allocation of page frames to processes
- Threshing, the working-set model
- Memory-mapped files and shared memory and
- Kernel memory allocation

Colorado State University

34

# Fritz-Rudolf Güntsch: Virtual Memory



Fritz-Rudolf Güntsch (1925-2012) at the Technische Universität Berlin in 1956 in his doctoral thesis, *Logical Design of a Digital Computer with Multiple Asynchronous Rotating Drums and Automatic High Speed Memory Operation*.

First used in Atlas, Manchester, 1962

PCs: Windows 95

When was Win 95 introduced?

Colorado State University

# 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 the same time
- Consider ability to execute partially-loaded program
  - Program no longer constrained by limits of physical memory
  - Each program uses less memory while running -> more programs run at the same time
    - Increased CPU utilization and throughput with no increase in response time or turnaround time
  - Less I/O needed to load or swap programs into memory -> each user program runs faster
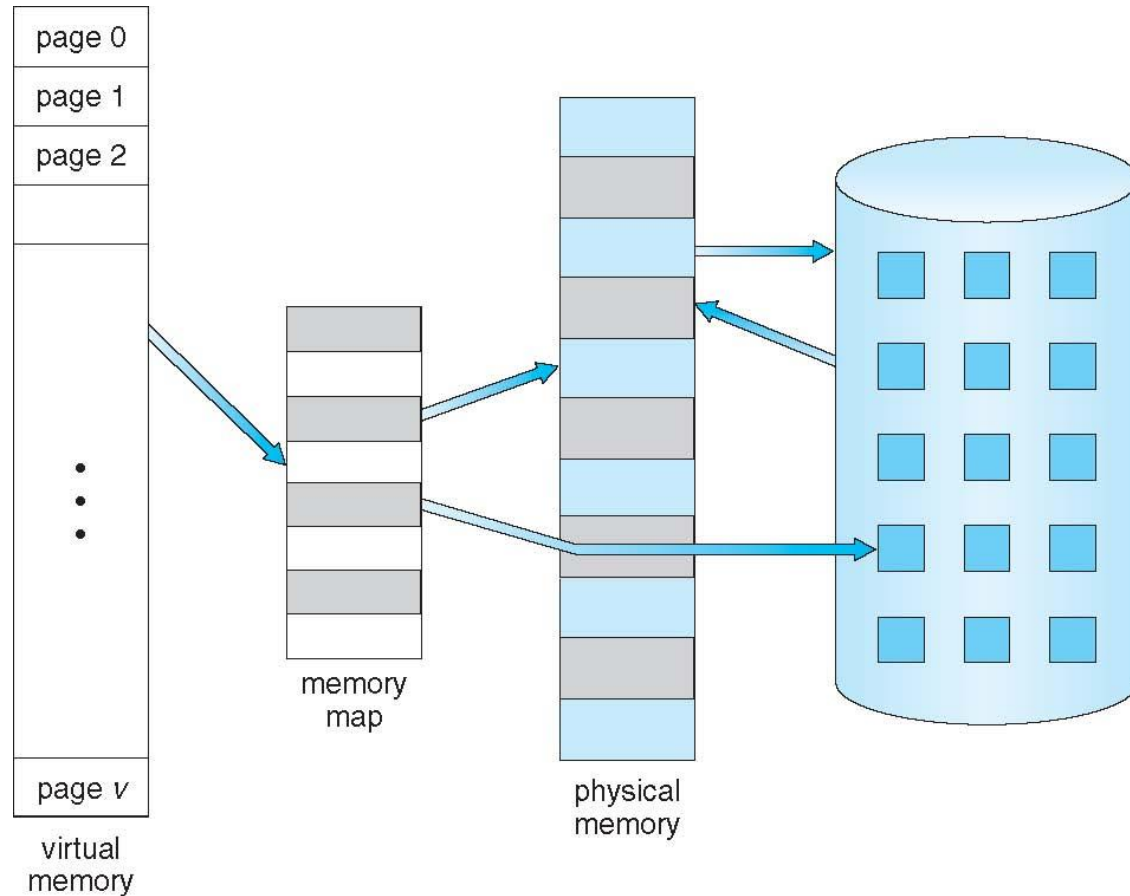
**Colorado State University**

# Background (Cont.)

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

- **Virtual address space** – logical view of how process views memory
  - Usually start at address 0, contiguous addresses until end of space
  - Meanwhile, physical memory organized in page frames
  - MMU must map logical to physical

- Virtual memory can be implemented via:
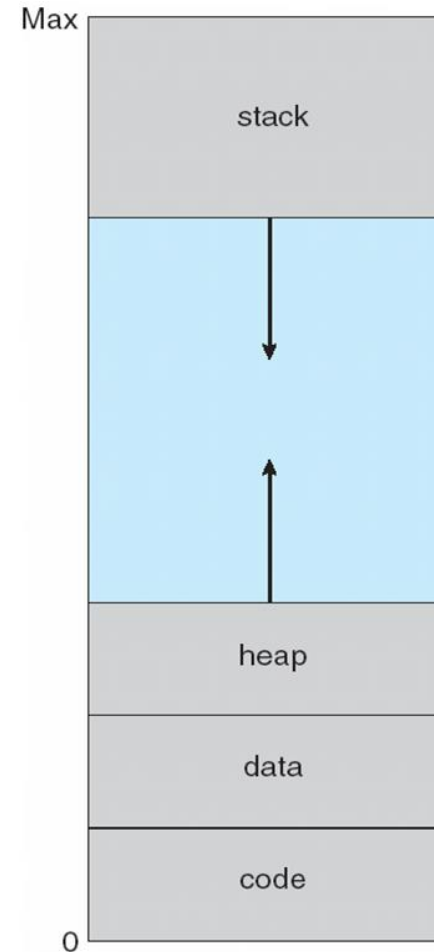  - Demand paging
  - Demand segmentation

That is the new idea

Colorado State University

page 0
page 1
page 2

page v

virtual memory

memory map

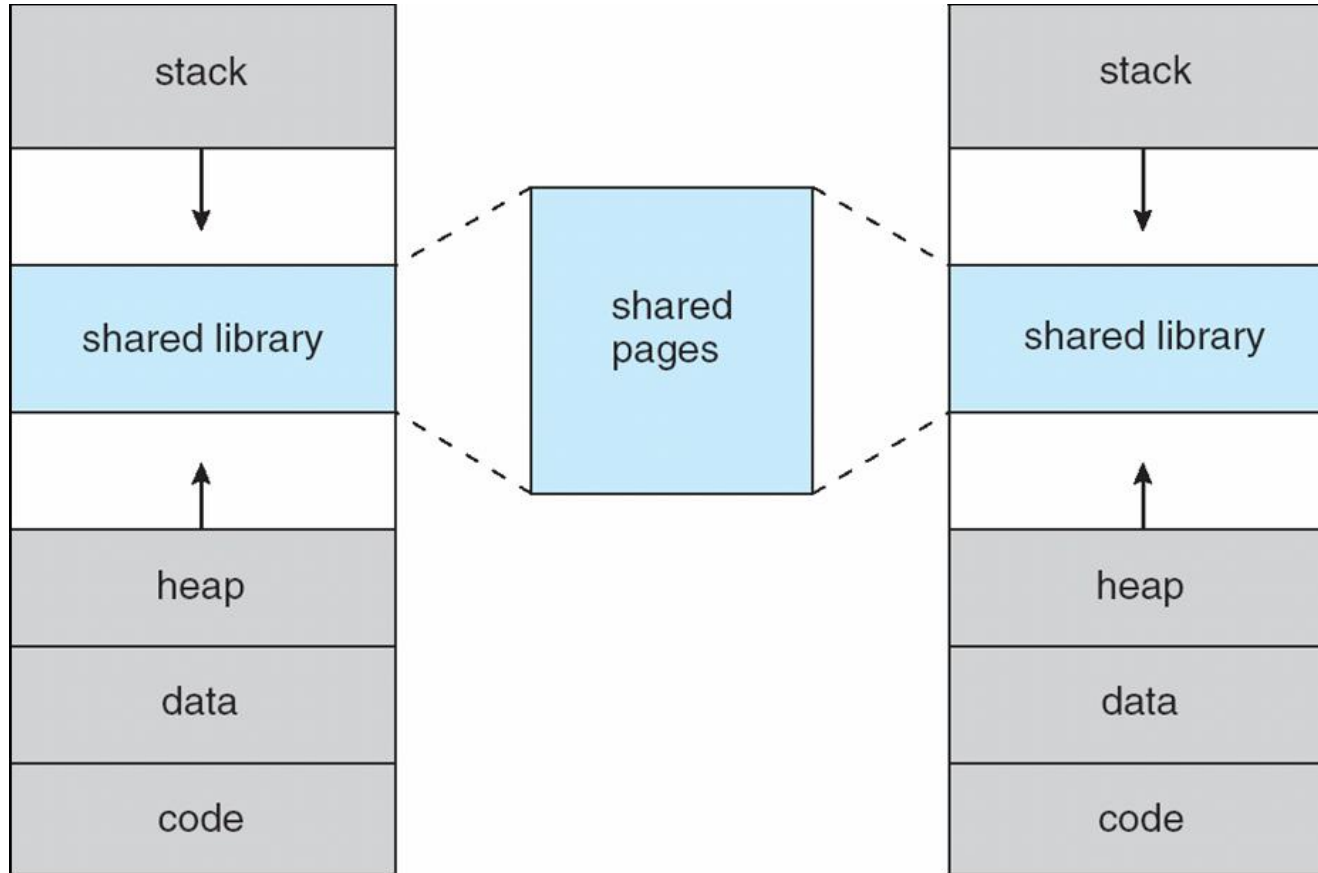physical memory

Colorado State University

# Virtual-address Space: advantages

- 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



**Colorado State University**

**Colorado State University**

# Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed: **Demand paging**
  - Less I/O needed, no unnecessary I/O
  - Less memory needed
  - Faster response
  - More users
- Similar to paging system with swapping
- 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**

**Colorado State University**

# Demand paging: Basic Concepts

- Demand paging: pager brings in only those pages into memory what are needed
- How to determine that set of pages?
  - Need new MMU functionality to implement demand paging
- 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

Colorado State University

# Valid-Invalid Bit

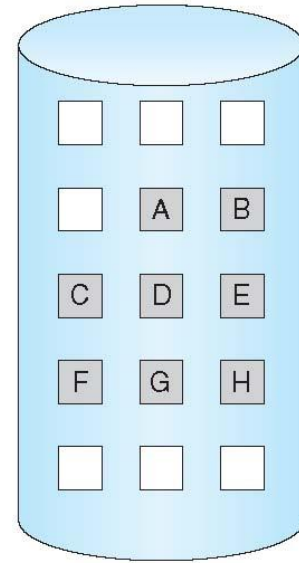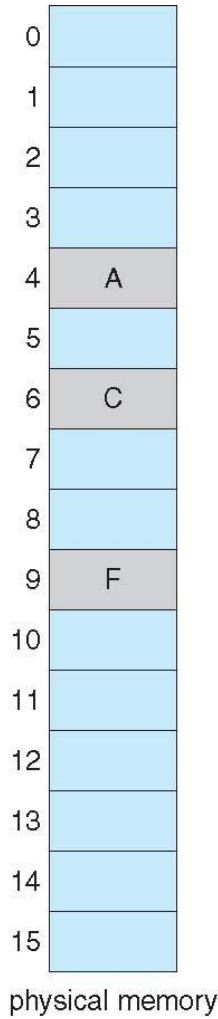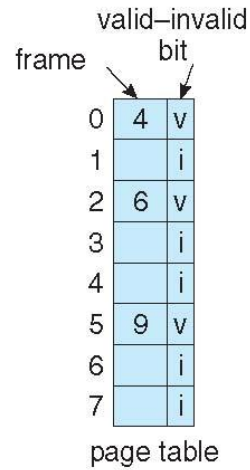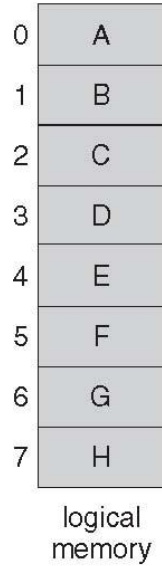- With each page table entry a valid–invalid bit is associated (**v** $\Rightarrow$ in-memory – **memory resident**, **i** $\Rightarrow$ not-in-memory)
- Initially valid–invalid bit is set to **i** on all entries
- Example of a page table snapshot:

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

page table

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

**Colorado State University**

Page 0 in Frame 4 (and disk)
Page 1 in Disk

Colorado State University

# Page Fault

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

**Page fault**

1. Operating system looks at a table to decide:
   - Invalid reference $\Rightarrow$ abort
   - Just not in memory, but in *backing storage*, ->2
2. Find free frame
3. Get 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

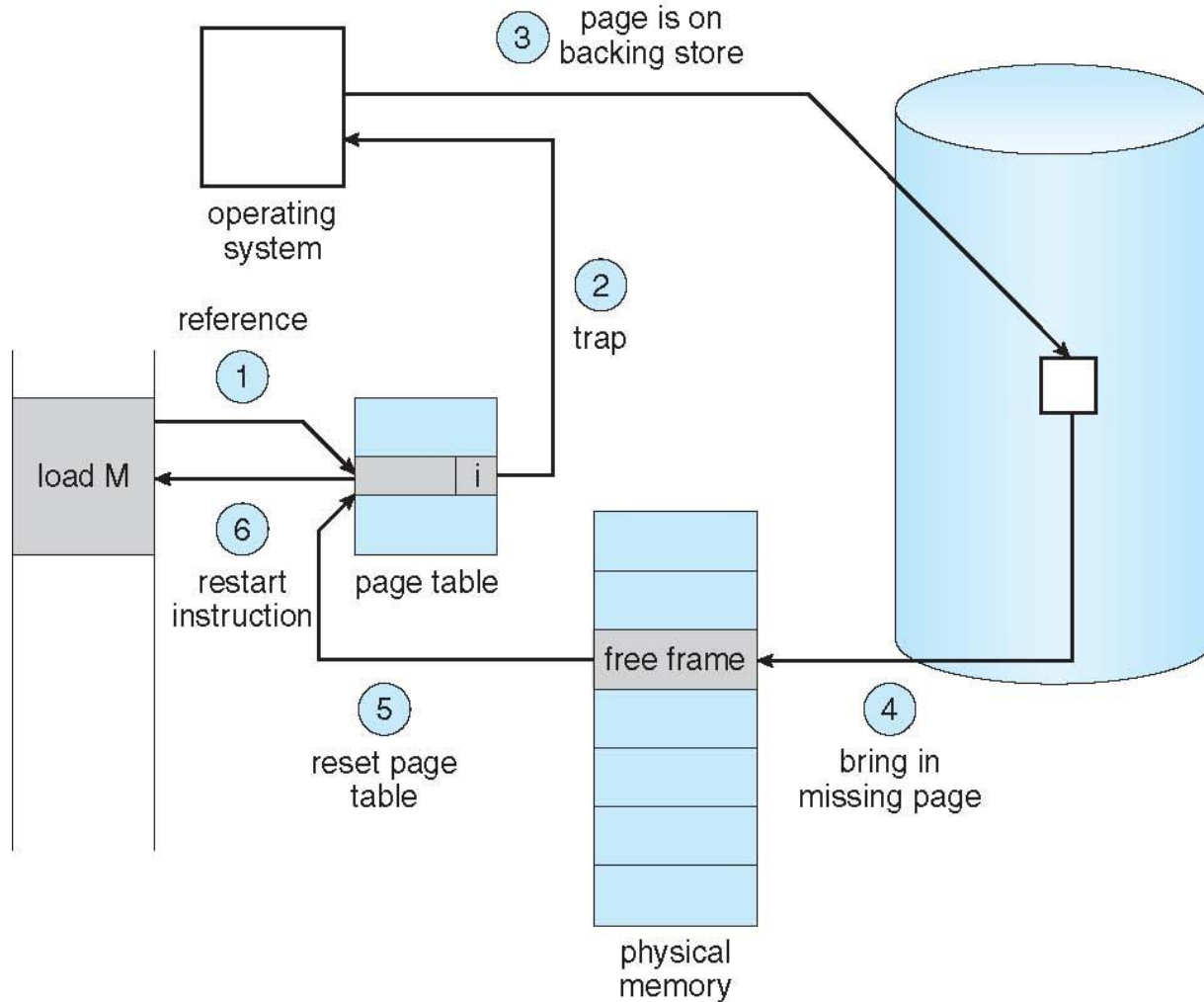Page fault: context switch because disk access is needed

Colorado State University

Solving a problem gives rise to a new class of problem:

- Contiguous allocation. Problem: external fragmentation

- Non-contiguous, but entire process in memory: Problem: Memory occupied by stuff needed only occasionally. Low degree of Multiprogramming.

- Demand Paging: Problem: page faults

- How to minimize page faults?

**Colorado State University**

Colorado State University

# Stages in Demand Paging (worse case)

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

- Three major activities
  - Service the interrupt – careful coding means just several hundred instructions needed
  - Read the page – relatively long time
  - Restart the process – again just a small amount of time
- Page Fault Rate $0 \leq p \leq 1$
  - if $p = 0$ no page faults
  - if $p = 1$, every reference is a fault
- Effective Access Time (EAT)

  EAT = $(1 - p)$ x memory access time

  $+ p$ (page fault overhead

  $+$ swap page out $+$ swap page in )

Hopefully p <<1

Page swap time = seek time + latency time

**Colorado State University**

# Demand Paging Simple Numerical Example

- Memory access time = 200 nanoseconds
- Average page-fault service time = 8 milliseconds
- EAT = (1 – p) x 200 ns + p (8 milliseconds)
  
      = (1 – p)  x 200 + p x 8,000,000  nanosec.
  
       = 200 + p x 7,999,800  ns

> Linear with page fault rate

- 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, **p = ?**
  - 220 > 200 + 7,999,800 x p
    20 > 7,999,800 x p
  - p < .0000025
  - < one page fault in every 400,000 memory accesses

We make some simplifying assumptions here.

**Colorado State University**