# Chapter 9: Virtual Memory

# Chapter 9:  Virtual Memory

- Background
- Demand Paging
- Copy-on-Write
- Page Replacement
- Allocation of Frames
- Thrashing
- Memory-Mapped Files
- Allocating Kernel Memory
- Other Considerations
- Operating-System Examples

# 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 -> more programs run at the same time

        - Increased CPU utilization and throughput with no increase in response time or turnaround time

# Background (Cont.)

- **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

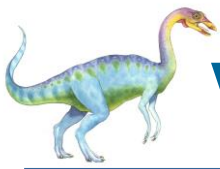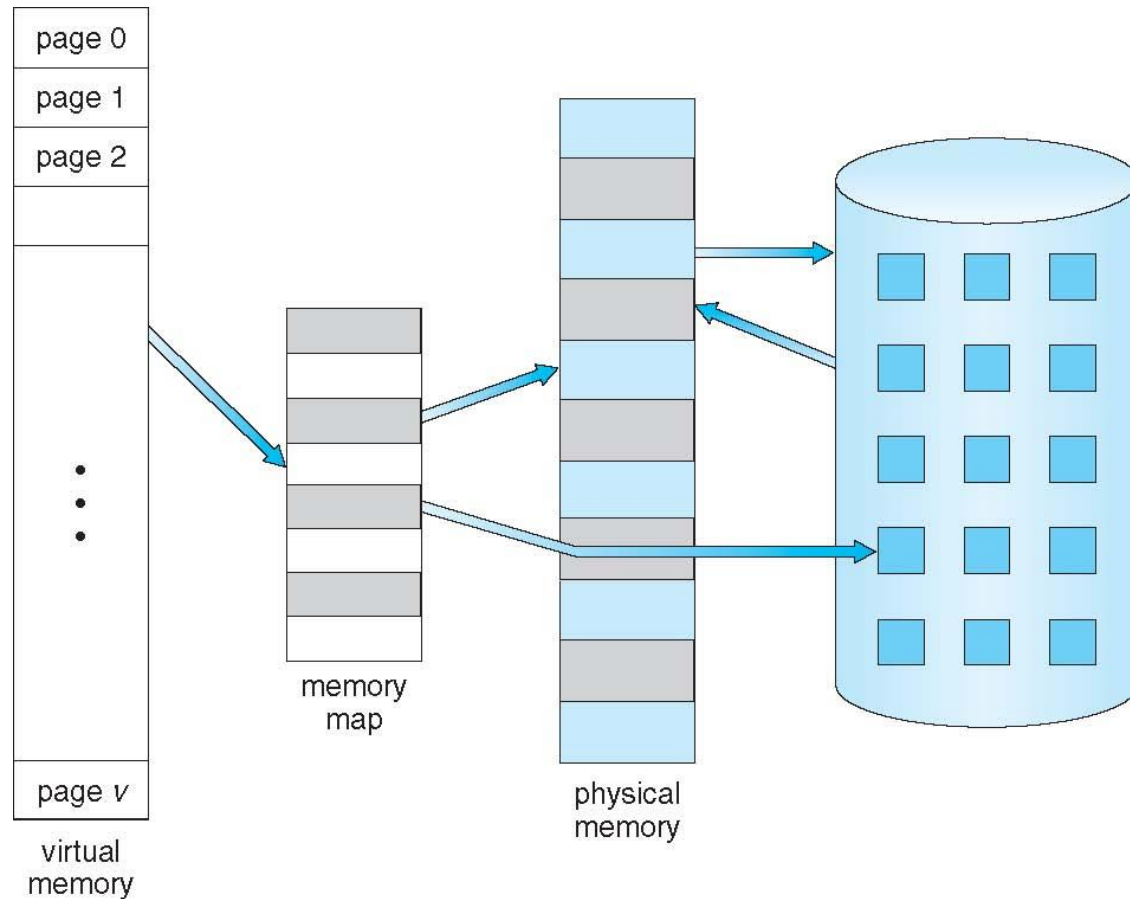# Background (Cont.)

- **Virtual address space** – logical view of how process is stored in memory

  - Usually starts 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:
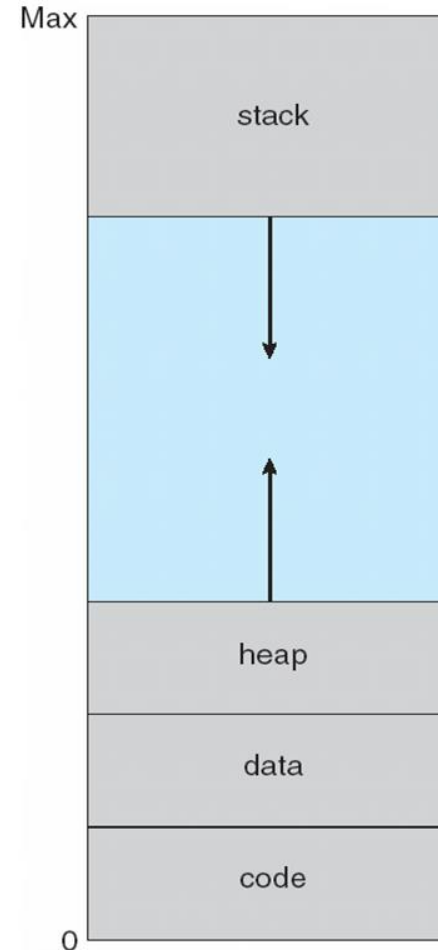
  - Demand paging

  - Demand segmentation

page 0
page 1
page 2

page v

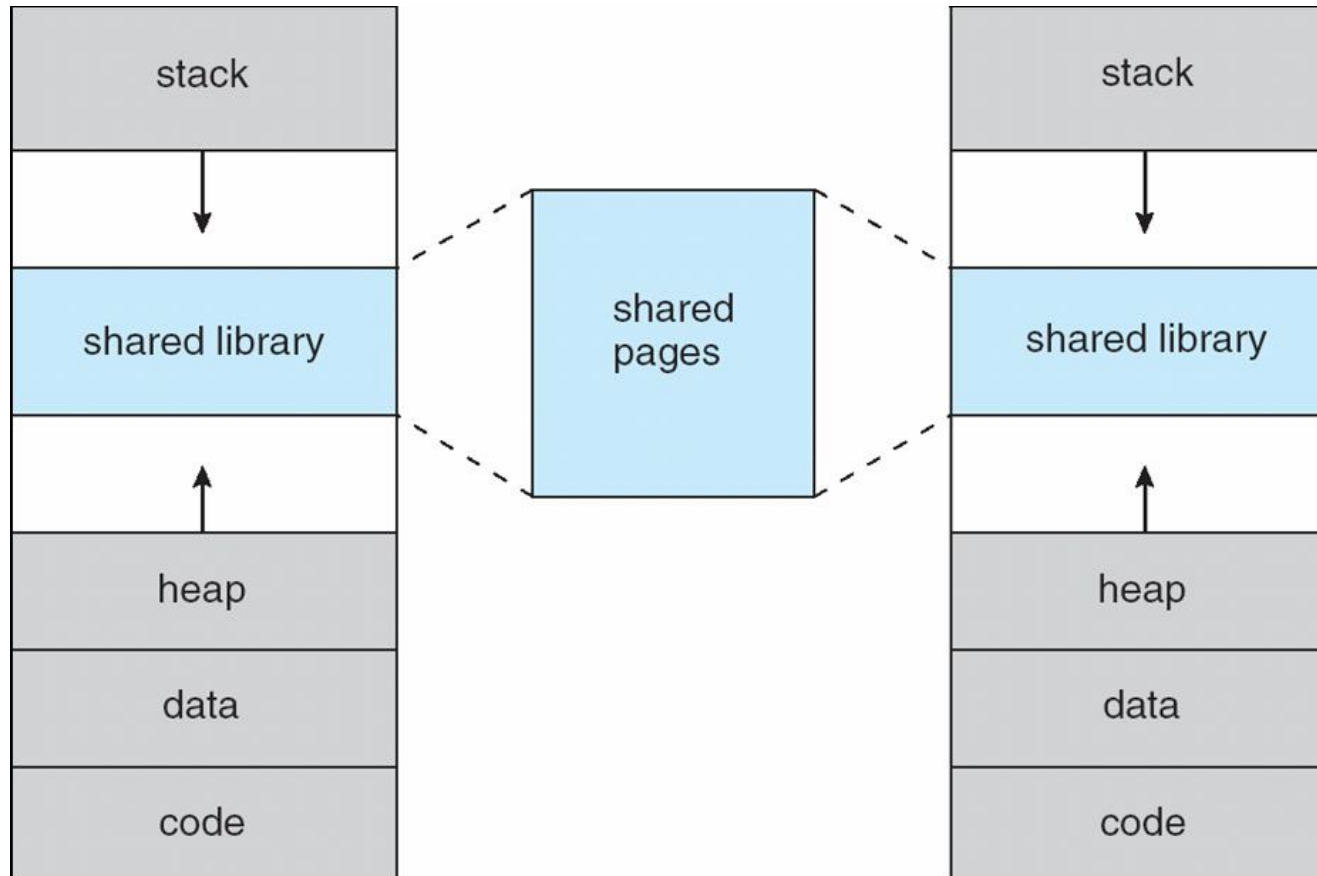virtual memory

memory map

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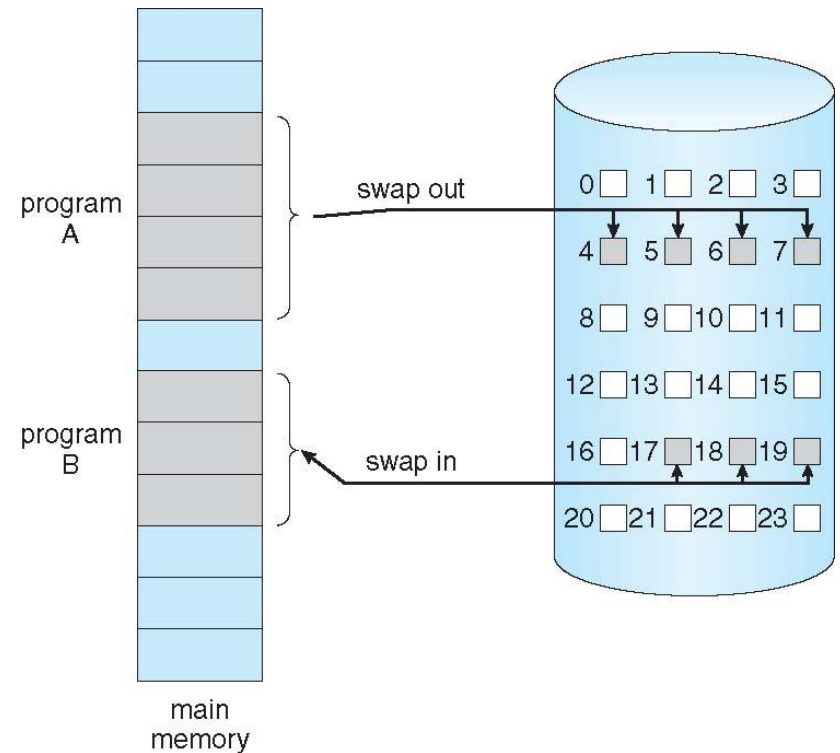- Pages can be shared during `fork()`, speeding process creation

# 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

- With swapping, pager guesses which pages will be used before swapping out again

- Instead, pager brings in only those pages into memory

- 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

# 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

# Page Fault

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

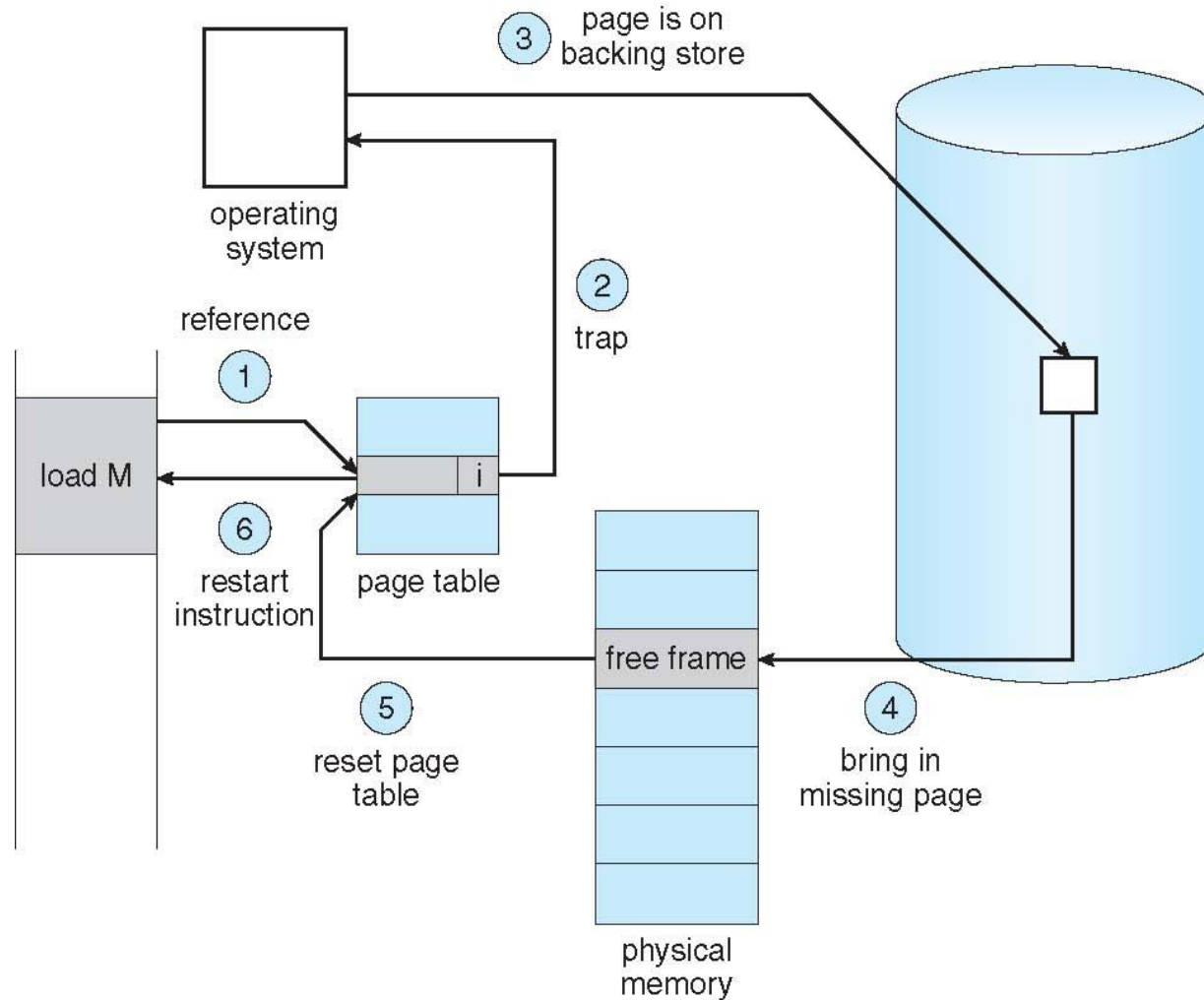    **page fault**

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

page is on backing store

3

operating system

reference

2

trap

1

load M

i

6

restart instruction

page table

free frame

5

reset page table

4

bring in missing page

physical memory

# 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
    - Consider fetch and decode of instruction which adds 2 numbers from memory and stores result back to memory
    - 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

# Performance of Demand Paging

- Stages in Demand Paging (worse case)
1. Trap to the operating system
2. Save 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. CPU is allocated to some other process
7. Receive an interrupt from the disk I/O subsystem (I/O completed)
8. Save the registers and process state for the other process
9. Determine that the interrupt was from the disk
10. Correct the page table and other tables to show that page is now in memory
11. Wait for the CPU to be allocated to this process again
12. Restore registers, process state, and new page table, and then resume the interrupted instruction

# Performance of Demand Paging (Cont.)

- Three major activities

  - Service the interrupt – careful coding means just several hundred instructions needed

  - Read the page – lots of 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)

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

# 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

# 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
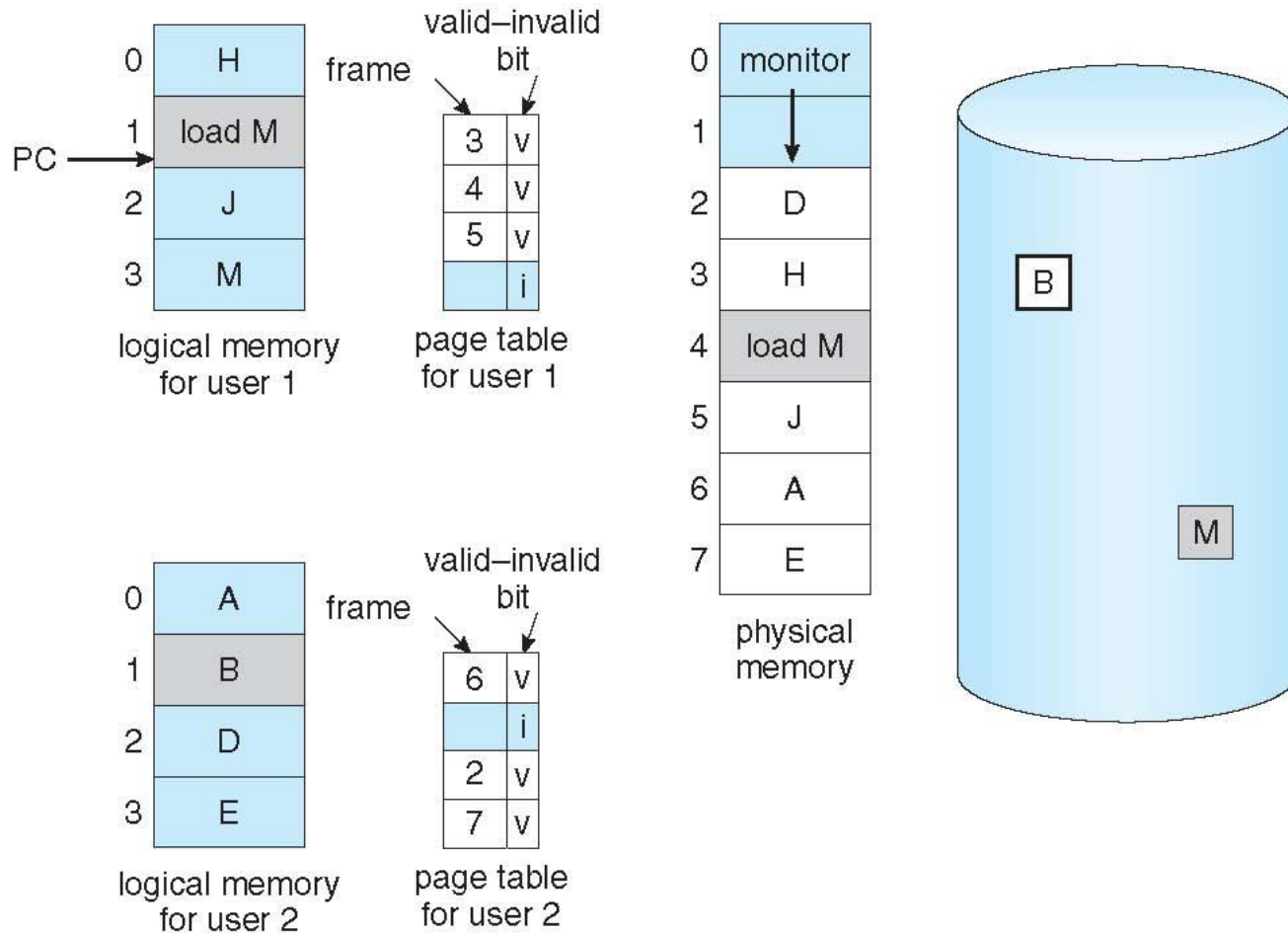
- Same page may be brought into memory several times

# 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



logical memory for user 1 · page table for user 1 · logical memory for user 2 · page table for user 2 · physical memory

# 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
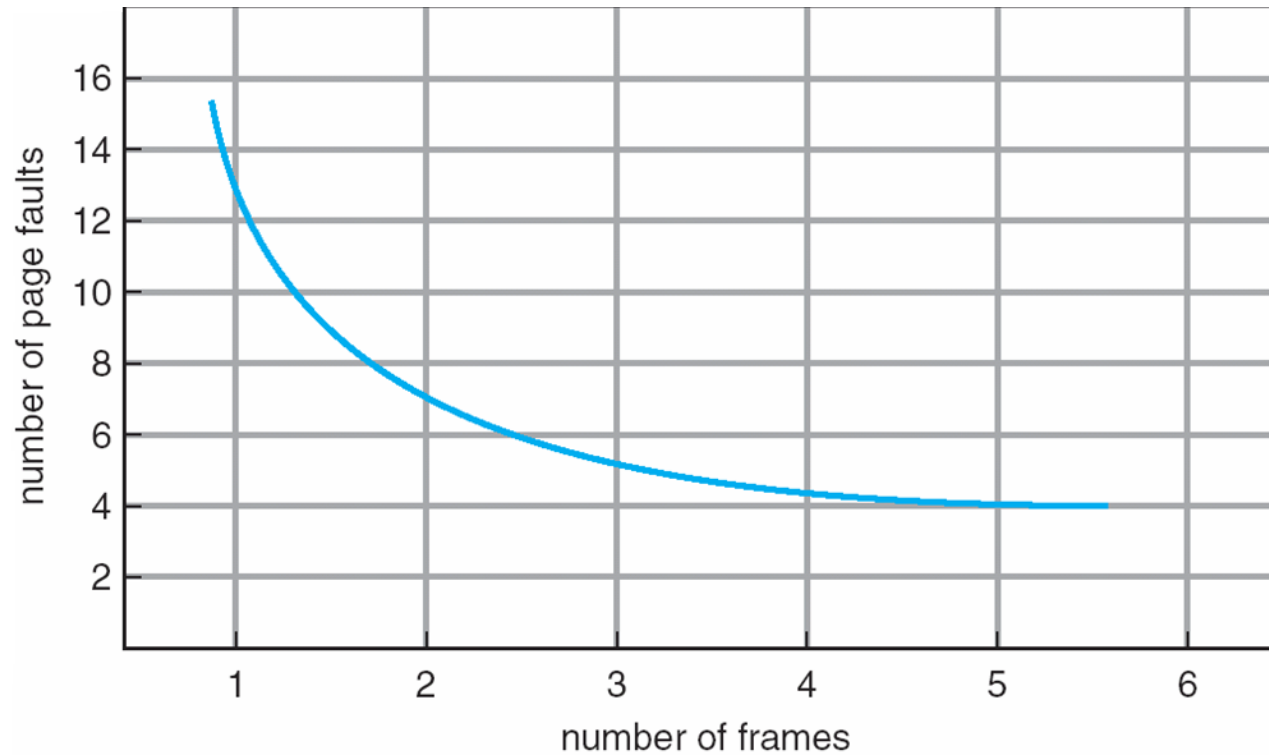
# Page and Frame Replacement Algorithms

- **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
  - Results depend on number of frames available
- In all our examples, the **reference string** of referenced page numbers is

<div align="center">

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

</div>

# 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)

reference string

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

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

page frames

15 page faults

# 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

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

page frames

- 12 faults – better than FIFO Generally good algorithm and frequently used

# Allocation of Frames

- Each process needs *minimum* number of frames
- *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 $p_i = \dfrac{s_i}{S} \times m$

$m = 64$

$s_1 = 10$

$s_2 = 127$

$a_1 = \dfrac{10}{137} \times 62 \approx 4$

$a_2 = \dfrac{127}{137} \times 62 \approx 57$

# 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
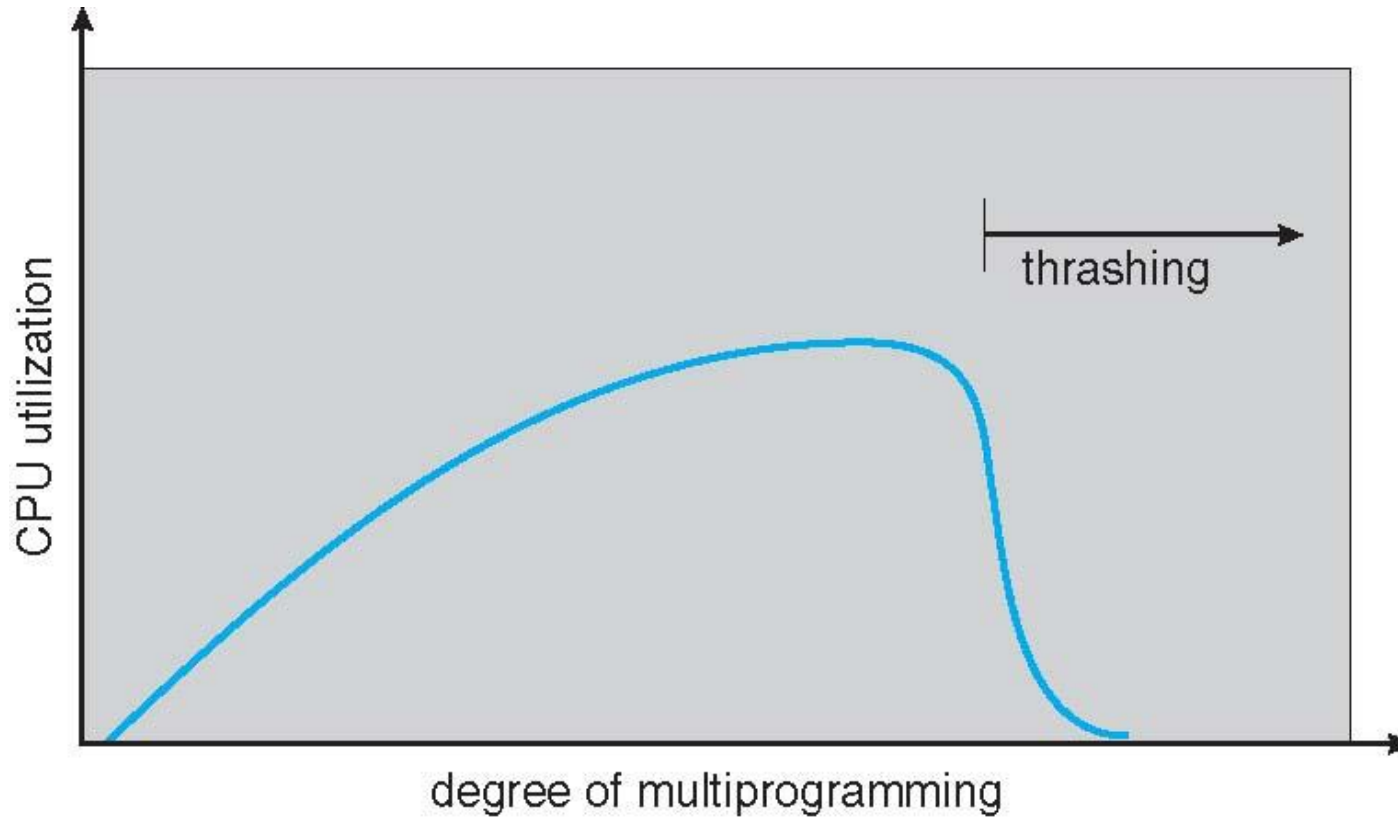
  - But possibly underutilized memory

# 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** $\equiv$ a process is busy swapping pages in and out

# Thrashing (Cont.)

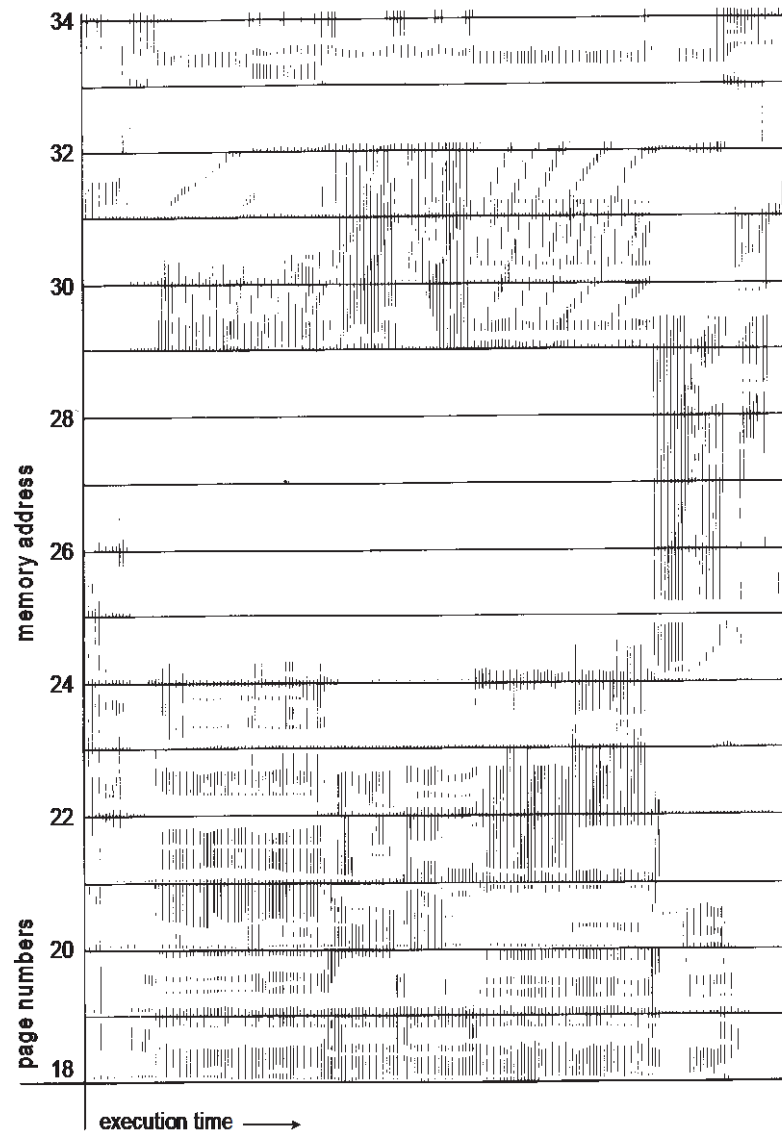# Demand Paging and Thrashing

- 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

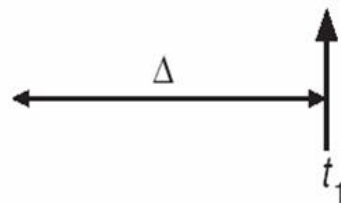  - Limit effects by using local or priority page replacement
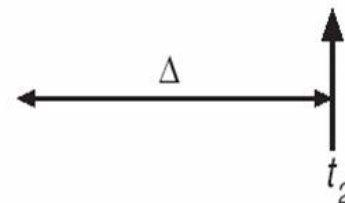
# Working-Set Model

- $\Delta \equiv$ 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

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 . . .

$\Delta$       $t_1$       $\Delta$       $t_2$

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

# Allocating Kernel Memory

- Treated differently from user memory

- Often allocated from a free-memory pool
    - Kernel requests memory for structures of varying sizes
    - Some kernel memory needs to be contiguous
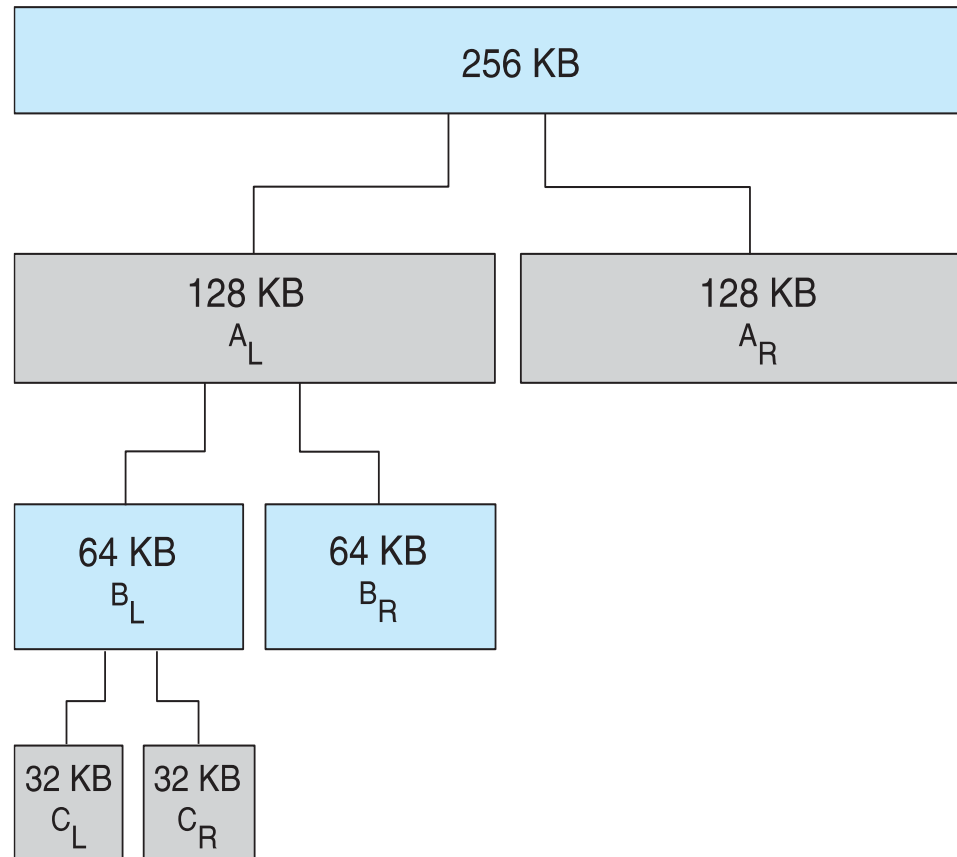        - I.e. for device I/O

# Buddy System

- Allocates memory from fixed-size segment consisting of physically-contiguous pages

- Memory allocated using **power-of-2 allocator**

  - Satisfies requests in units sized as power of 2

  - Request rounded up to next highest power of 2

  - When smaller allocation needed than is available, current chunk split into two buddies of next-lower power of 2

    ▸ Continue until appropriate sized chunk available

- For example, assume 256KB chunk available, kernel requests 21KB

  - Split into $A_L$ and $A_R$ of 128KB each

    ▸ One further divided into $B_L$ and $B_R$ of 64KB

      – One further into $C_L$ and $C_R$ of 32KB each – one used to satisfy request

- Advantage – quickly **coalesce** unused chunks into larger chunk

- Disadvantage - fragmentation
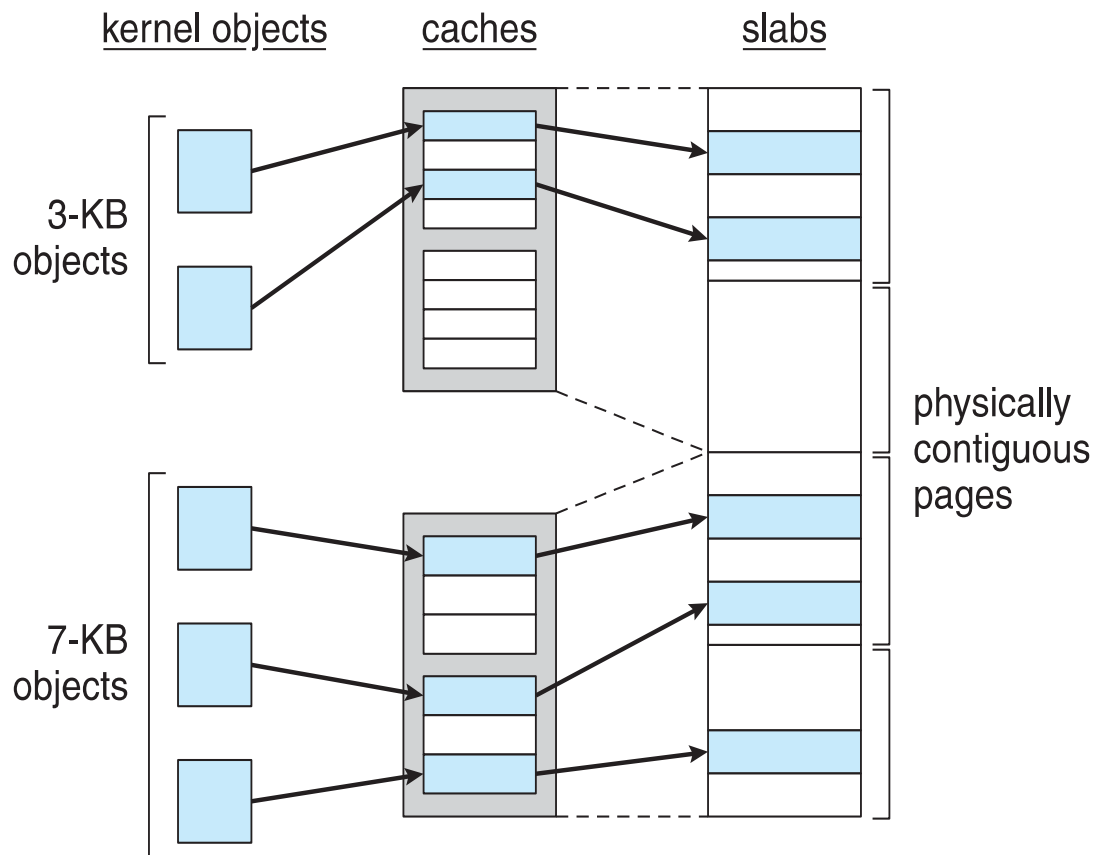
# Buddy System Allocator

physically contiguous pages

# Slab Allocator

- Alternate strategy

- **Slab** is one or more physically contiguous pages

- **Cache** consists of one or more slabs

- Single cache for each unique kernel data structure

  - Each cache filled with **objects** – instantiations of the data structure

- When cache created, filled with objects marked as `free`

- When structures stored, objects marked as `used`

- If slab is full of used objects, next object allocated from empty slab

  - If no empty slabs, new slab allocated

- Benefits include no fragmentation, fast memory request satisfaction

# Slab Allocation



kernel objects    caches    slabs

3-KB objects

7-KB objects

physically contiguous pages

# Slab Allocator in Linux

- For example process descriptor is of type `struct task_struct`
- Approx 1.7KB of memory
- New task -> allocate new struct from cache
    - Will use existing free `struct task_struct`
- Slab can be in three possible states
    1. Full – all used
    2. Empty – all free
    3. Partial – mix of free and used
- Upon request, slab allocator
    1. Uses free struct in partial slab
    2. If none, takes one from empty slab
    3. If no empty slab, create new empty

# Slab Allocator in Linux (Cont.)

- Slab started in Solaris, now wide-spread for both kernel mode and user memory in various OSes

- Linux 2.2 had SLAB, now has both SLOB and SLUB allocators

  - SLOB for systems with limited memory

    - Simple List of Blocks – maintains 3 list objects for small, medium, large objects

  - SLUB is performance-optimized SLAB removes per-CPU queues, metadata stored in page structure

# End of Chapter 9