# Virtual Memory

---

# Background

- Code needs to be in memory to execute, but entire program rarely used
  - 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
  - Less I/O needed to load or swap programs into memory -> each user program runs faster
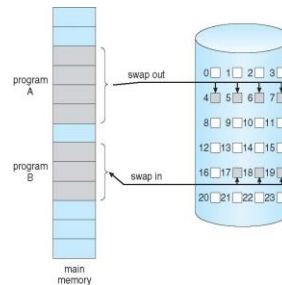
---

# Virtual memory

- **Virtual 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
  - More programs running concurrently
  - Less I/O needed to load or swap processes

# Virtual memory  (Cont.)

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

# 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 $\Rightarrow$ 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**

# 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
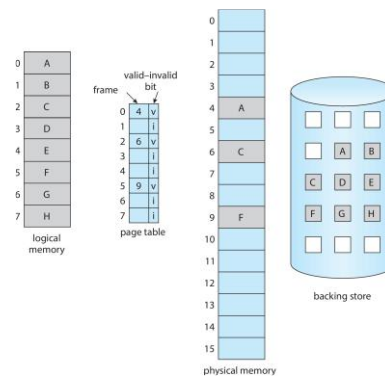    - Without programmer needing to change code

# 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:



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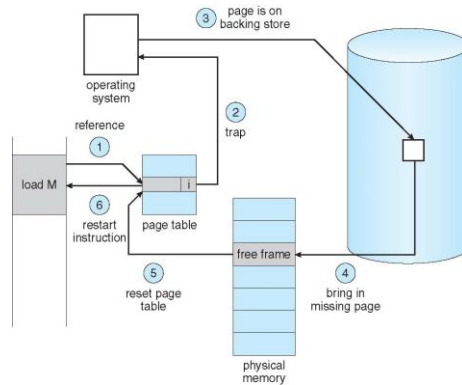
---

# Page Table When Some Pages Are Not in Main Memory

---

# Steps in Handling Page Fault

1. If there is a reference to a page, first reference to that page will trap to operating system
   - *Page fault*
2. Operating system looks at another table to decide:
   - Invalid reference $\Rightarrow$ abort
   - Just not in memory
3. Find free frame
4. Swap page into frame via scheduled disk operation
5. Reset tables to indicate page now in memory
   Set validation bit = **v**
6. Restart the instruction that caused the page fault

# Steps in Handling a Page Fault (Cont.)

---

# 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
- Hardware support needed for demand paging
  - Page table with valid / invalid bit
  - Secondary memory (swap device with **swap space**)
  - Instruction restart

---

# Free-Frame List

- When a page fault occurs, the operating system must bring the desired page from secondary storage into main memory.
- Most operating systems maintain a **free-frame list** -- a pool of free frames for satisfying such requests.

head $\longrightarrow$ 7 $\longrightarrow$ 97 $\longrightarrow$ 15 $\longrightarrow$ 126 ... $\longrightarrow$ 75

- Operating system typically allocate free frames using a technique known as **zero-fill-on-demand** --  the content of the frames zeroed-out before being allocated.
- When a system starts up, all available memory is placed on the free-frame list.

## 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

## Stages in Demand Paging  (Cont.)

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

## Performance of Demand Paging

- 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

    $+ p$ (page fault overhead

    + swap page out

    + 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)

$$= (1 – p) \text{ x } 200 + p \text{ x } 8{,}000{,}000 \text{ (nanosec)}$$
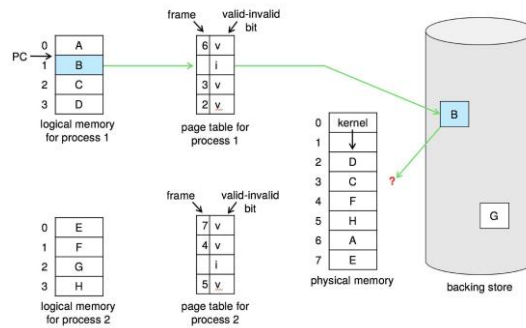$$= 200 + p \text{ x } 7{,}999{,}800 \text{ (nanosec)}$$

---

# What Happens if There is no Free Frame?

- Used up by process pages

- 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
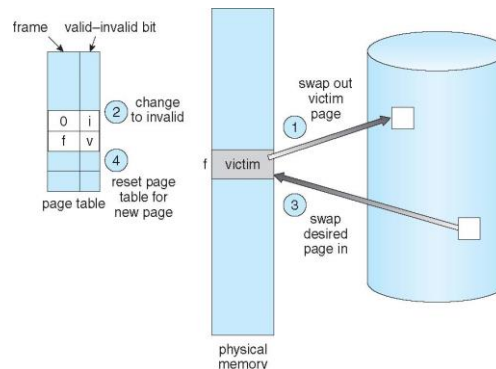
# Need For Page Replacement

---

# 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

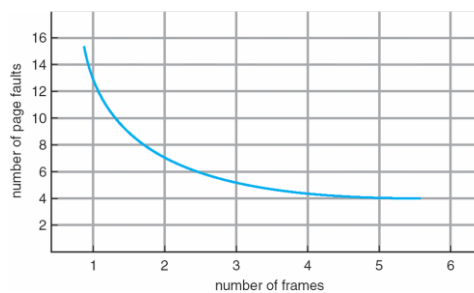Note now potentially 2 page transfers for page fault – increasing EAT

---

# 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
  - ☐ Results depend on number of frames available
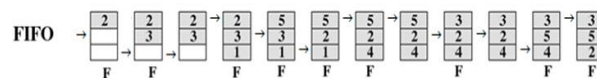
- ☐ Page-replacement algorithms: OPT(MIN), LRU, FIFO, CLOCK

---

## Graph of Page Faults Versus The Number of Frames
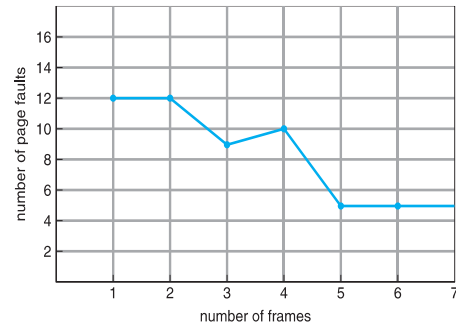
---

## *FIFO* algorithm

- ☐ View frames allocated to the process as a circular buffer
  - ☐ When the buffer is full, the oldest page is replaced: first-in first-out
  - ☐ Simple implementation: just a pointer to cycle the frames
- ☐ Reference string:      2 3 2 1 5 2 4 4 5 3 2 5 2
- ☐ 3 frames (3 pages can be in memory at a time per process)



- ☐ Adding more frames can cause more page faults!
  - ☐ **Belady's Anomaly**

# FIFO Illustrating Belady's Anomaly

---

# *OPT* algorithm

- The page that will not be referenced for the longest time is replaced

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

- Used for measuring how well your algorithm performs

---

# *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



- 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 Most Recent Page References

reference string

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

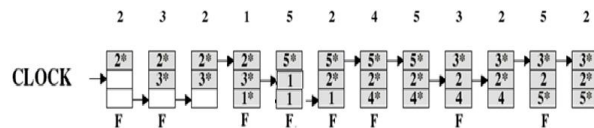| stack before a | stack after b |
|---|---|
| 2 | 7 |
| 1 | 2 |
| 0 | 1 |
| 7 | 0 |
| 4 | 4 |

a   b

---

# LRU Approximation Algorithms

- LRU needs special hardware and still slow
- **Reference bit / use 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) algorithm

- The frames for the process are treated as a circular buffer.

- When a page is changed, the pointer will point to the next frame in the buffer.

- Each frame has a *use* bit. This bit is set to 1 when
  - A page is first loaded into the frame
  - The page in the frame is referenced

- The page is only replaced at the frame with *use bit = 0*
  - While trying to find the page to replace, all *use bits* are reset to 0

## second chance (CLOCK) algorithm



NOTE: Asterisk (*) indicates *reference/use bit = 1*

## Enhanced Second-Chance Algorithm

- Improve algorithm by using reference bit and modify bit (if available) in concert
- Take ordered pair (reference, modify):
  - (0, 0) neither recently used not modified – best page to replace
  - (0, 1) not recently used but modified – not quite as good, must write out before replacement
  - (1, 0) recently used but clean – probably will be used again soon
  - (1, 1) recently used and modified – probably will be used again soon and need to write out before replacement
- When page replacement called for, use the clock scheme  but use the four classes replace page in lowest non-empty class
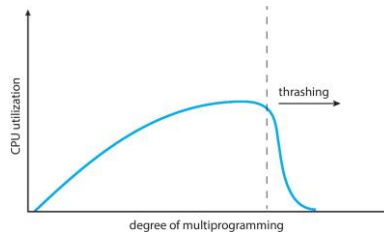  - Might need to search circular queue several times

# 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 (Cont.)

- **Thrashing**. A process is busy swapping pages in and out

# 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

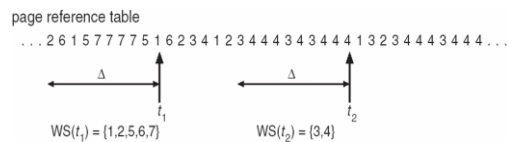- Limit effects by using local or priority page replacement

# 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

# Working-Set Model (Cont.)

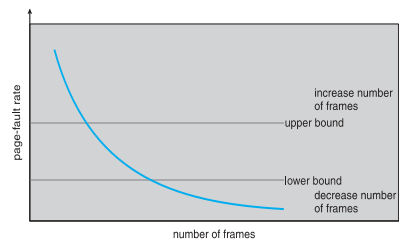- if $D > m \Rightarrow$ Thrashing
- Policy if $D > m$, then suspend or swap out one of the processes

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 . . .

$\Delta$      $\Delta$

$t_1$      $t_2$

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

# Page-Fault Frequency

- More direct approach than WSS
- Establish "acceptable" **page-fault frequency** (**PFF**) rate and use local replacement policy
  - If actual rate too low, process loses frame
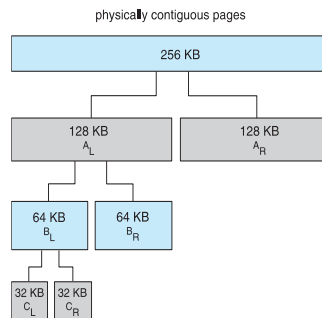  - If actual rate too high, process gains frame

increase number of frames

upper bound

lower bound

decrease number of frames

page-fault rate

number of frames

# 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