# 10. Virtual Memory

ECE30021/ITP30002 Operating Systems

# Agenda

- Background
- Demand paging
- Copy-on-Write (COW)
- Page replacement
- Allocation of frames
- Thrashing
- Memory-mapped files
- Allocating kernel memory
- Other considerations
- Operation-system examples

# Background

- **Instructions should be in physical memory to be executed**
  - ➔ In order to execute a program, should we load <span style="color:red">entire program</span> in memory?

- **Some parts are rarely used**
  - Error handling codes
  - Arrays/lists larger than necessary
  - Rarely used routines

- **Alternative**
  - <span style="color:red">Executing program which is only partially in memory</span>
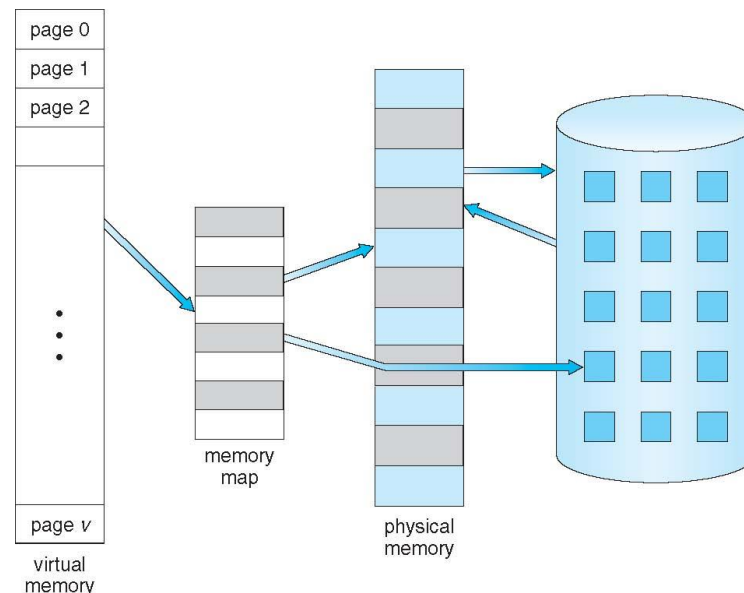
# Background

- **If we can run a program by loading in parts …**
  - A program is not constrained by the amount of physical memory
  - More program can run at the same time
  - Less I/O is need to load or swap programs

# Virtual Memory

- **Virtual memory**: a separation of logical memory from physical memory
  - Contents not currently reside in main memory can be addressed.
  - ➔ H/W and OS will load the required memory from auxiliary storage automatically.
  - User programs can reference more memory than actually exists



page 0
page 1
page 2

page v

virtual memory
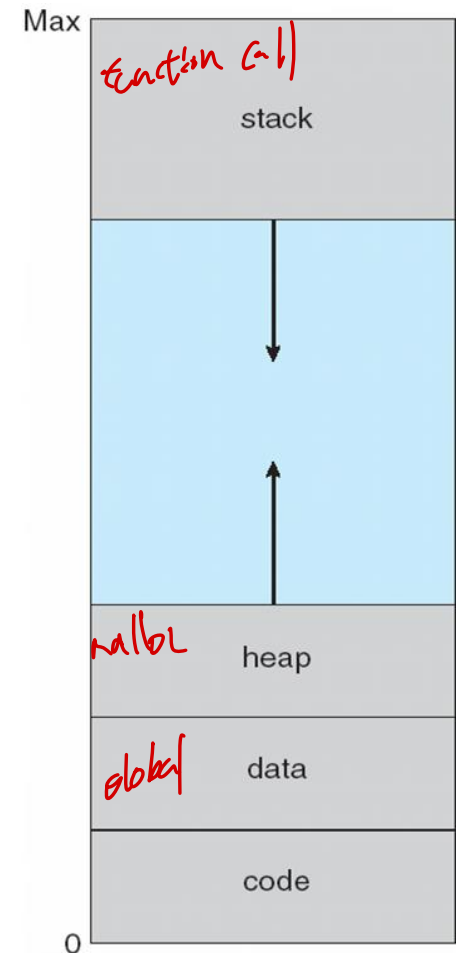
memory map

physical memory

# Virtual Memory

- **Virtual memory**: a separation of logical memory from physical memory
  - Only part of the program needs to be in memory for execution
  - Logical address space can be much larger than physical address space
  - Allows address spaces to be shared by several processes
  - Allows for more efficient process creation (e.g. COW)
  - More programs running concurrently
  - Less I/O needed to load or swap processes
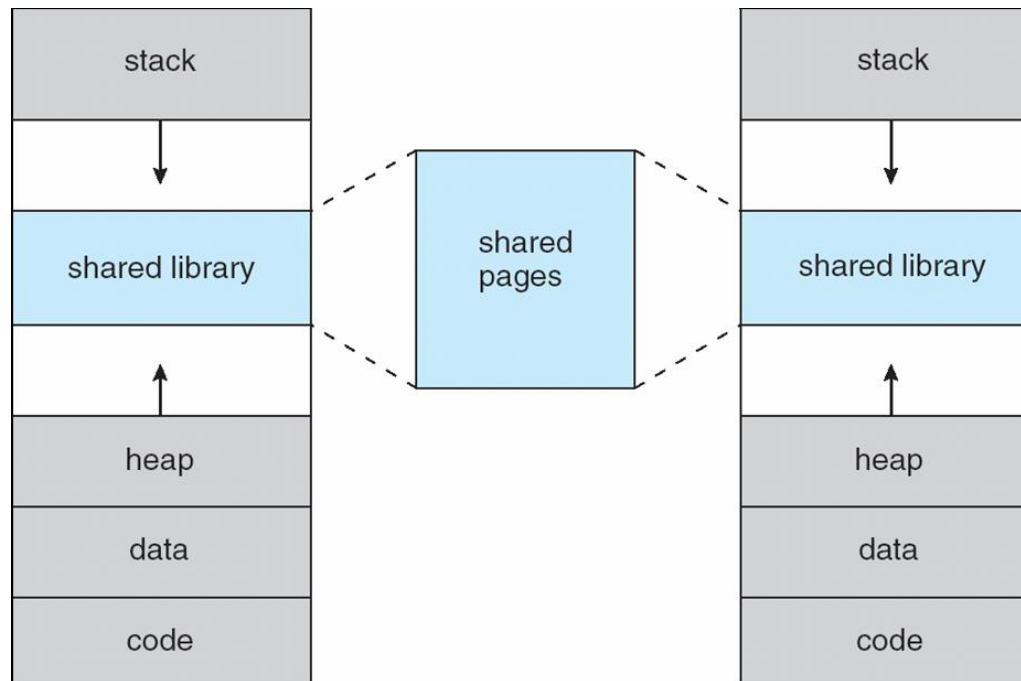
    page 로드

# Virtual Address Space

■ **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
- Programmers don't have to concern about memory management
- Virtual address space can be sparse

Max

*function call*

stack

*malloc*

heap

*global*

data

code

0

# Shared Library Using Virtual Memory

- Shared page

# Implementation Techniques

■ **Demand paging**
  ■ Paging + swapping

■ **Demand segmentation**
  ■ Segmentation + swapping

page 0
page 1
page 2

⋮

page *v*

virtual
memory

memory
map
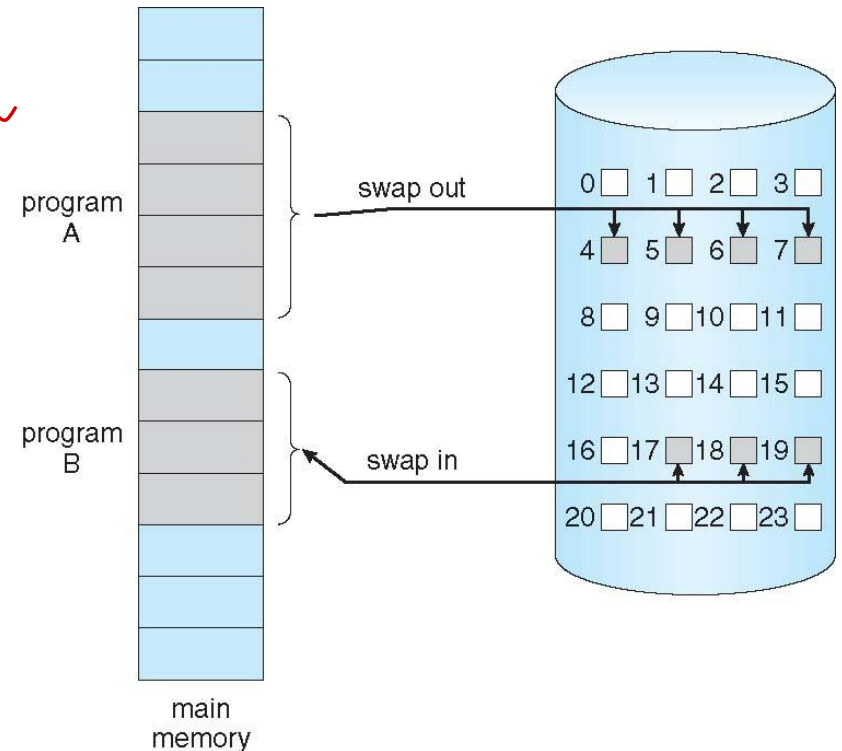
physical
memory

block 크기 = page 크기

# Agenda

- Background
- <u>Demand paging</u>
- Copy-on-Write (COW)
- Page replacement
- Allocation of frames
- Thrashing
- Memory-mapped files
- Allocating kernel memory
- Other considerations
- Operation-system examples

# Demand Paging

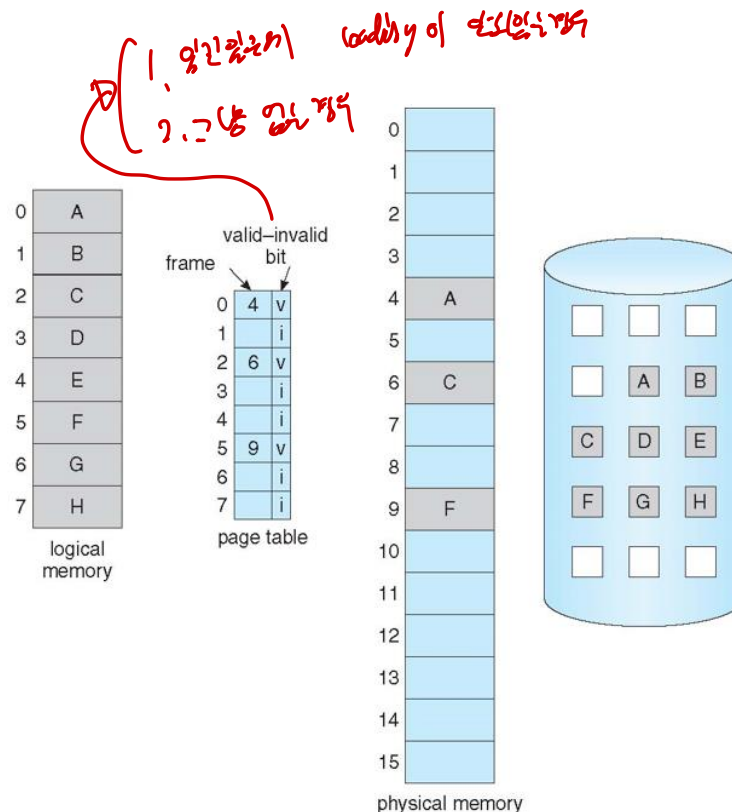- **Similar to paging system with swapping**

꼭 필요할때만 가져나

page swapper

- **Lazy swapper (or pager)**
  - Pages are only loaded when they are demanded during execution.

# Demand Paging

- **Pages are loaded only when they are demanded during program execution**
  - Requires H/W support to distinguish the page on memory or on disk



logical memory · page table · physical memory

# Page Table with Valid/Invalid Bit

- **Valid/invalid bit of each page**
    - **Valid**: the page is valid and exists in physical memory
    - **Invalid**: the page is not valid (not in the valid logical address space of the process) or not loaded in physical memory

- **If program tries to access ..**
    - Valid page: execution proceeds normally
    - Invalid page: cause **page-fault** trap to OS
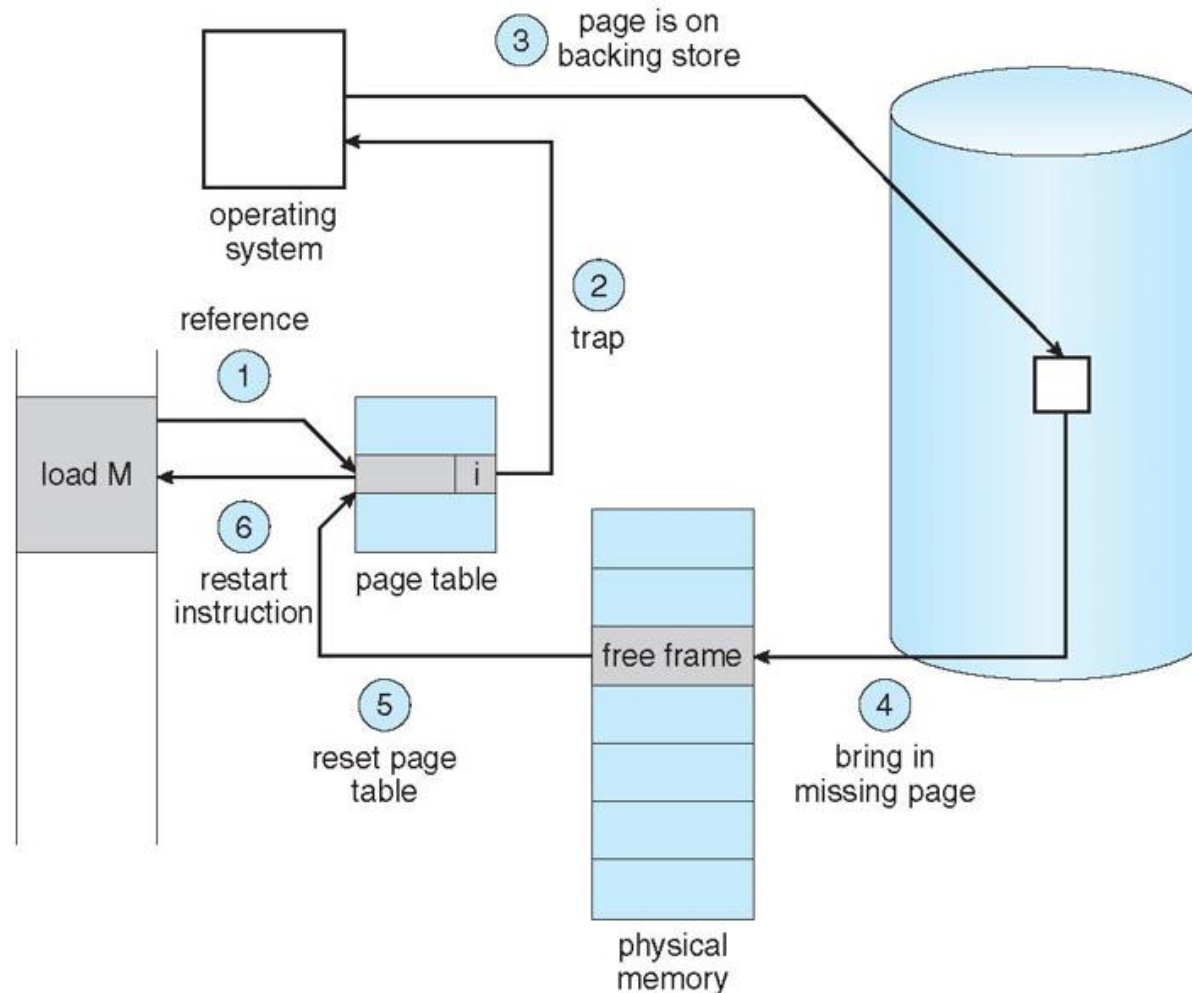
Frame #        valid-invalid bit

| | |
|---|---|
| | |
| | v |
| | v |
| | v |
| | i |
| . . . | |
| | i |
| | i |

page table

# Basic Concepts

page fault 발생시 valid-invalid bit이 | 경우 참조한다. ⌈ 종료
⌊ 로드

- **Handling page-fault**
  - Check an internal table to determine whether the reference was valid or not
  - If the reference was invalid, terminates the process
  - If valid, page it in.
    - □ Find a free frame
    - □ Read desired page into the free frame
    - □ Modify internal table
    - □ Restart the instruction that caused the page-fault trap

# Handling Page−Fault

# Basic Concepts

- **In pure demand paging, a page is not brought into memory until it is required**

- **Performance degradation**
  - Theoretically, some program can cause multiple page faults per instructions
  - But actually, this behavior is exceedingly unlikely.
    - locality of references

- **H/W supports for demand paging**
  - Page table (support for valid-invalid bit, …)
  - Secondary memory (swap space)
  - Instruction restart

# Performance of Demand Paging

- **Effective access time**

  Effective access time = $(1-p)*ma + p * $ <page fault time>

  - *ma* : ^main^ memory access time (10~200 nano sec.)
  - *p* : probability of page fault ~~ex~~) 0.00

- **Page fault time**
  - Service page-fault interrupt      -> 1~100 μsec.
  - Read in the page      -> about 8 msecs
  - Restart the process      -> 1~100 μsec.

# Performance of Demand Paging

Effective access time = $(1-p) * ma + p * $ <page fault time>

- **Example**
  - Memory access time: 200 nano sec
  - Page-fault service time: 8 milliseconds

- **Then···**

  Effective access time (in nano sec.)

  $= (1-p) * 200 + 8{,}000{,}000 * p$

  $\approx 200 + 7{,}999{,}800 * p$

  - Proportional to **page fault rate**

    M a x u D

    Ex) p == 1/1000, effective access time = 8.2 μsec. (40 times)

  - Page fault rate should be kept low

# Execution of Program in File System

- **Ways to execute a program in file system**
    - Option1: copy entire file into swap space at starting time
        - Usually swap space is faster than file system
    - Option2: initially, demand pages from files system and all subsequent paging can be done from swap space
        - Only needed pages are read from file system

# Agenda

- Background
- Demand paging
- <u>Copy-on-Write (COW)</u>
- Page replacement
- Allocation of frames
- Thrashing
- Memory-mapped files
- Allocating kernel memory
- Other considerations
- Operation-system examples

# Copy-on-Write

- **fork() copies process**
  - Duplicates pages belong to the parent

*process creation, 마르지 발생 overhead 크다* *parent process가 child process로 복사되면서 한꺼번에 복사되다 원래라면 좋은 내용 저장 되면서*
*특히 그냥 process가 writing를 하면 데이터 내용 수정하며 한다 복제할 수 없고 실제 복제되면 복제되며 memory content 발생 -한꺼번에다.*

- **Copy-on-write (COW)**
  - When the process is created, pages are not actually duplicated but just shared.
    - Process creation time is reduced.
  - When either process writes to a shared page, a copy of the page is created.

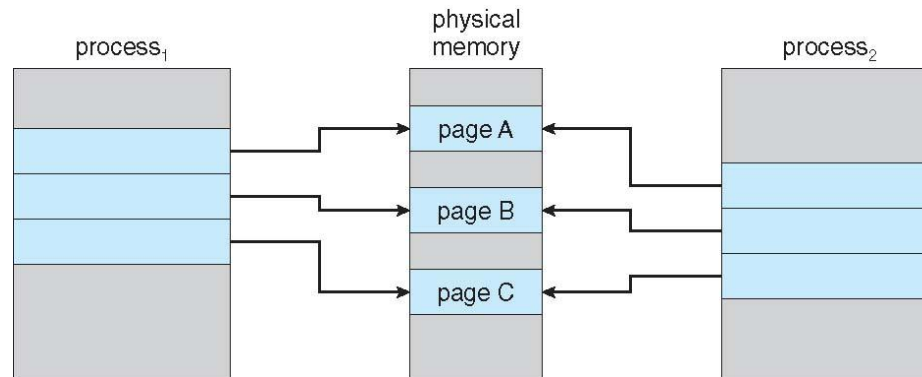  cf. vfork() – logically shares memory with parent (<u>obsolete</u>)

- **Many OS's provides a list of free frames for COW or stack/heap that can be expanded**
  - ➔ Zero-fill-on-demand (ZFOD)
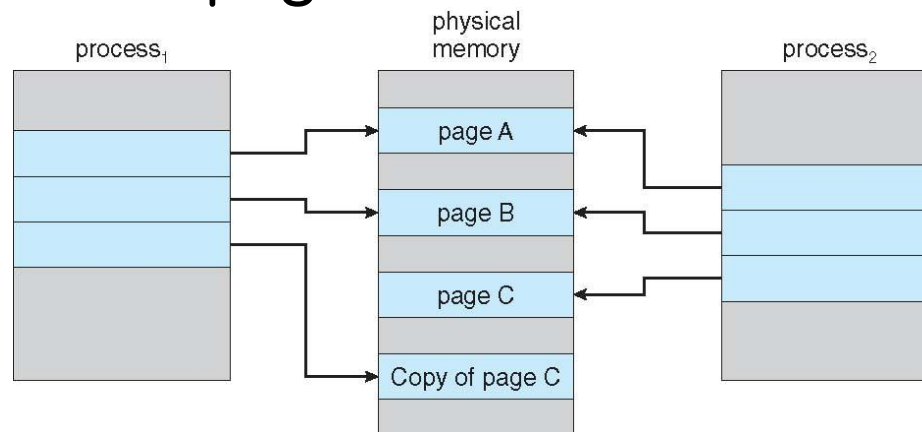    - Zero-out pages before being assigned to a process

# Copy-on-Write

■ Before P1 modifies page C



■ After P1 modifies page C

# Agenda

- Background
- Demand paging *paging + swapping*
- Copy-on-Write (COW)
- <u>Page replacement</u>
- <u>Allocation of frames</u>
- Thrashing
- Memory-mapped files
- Allocating kernel memory
- Other considerations
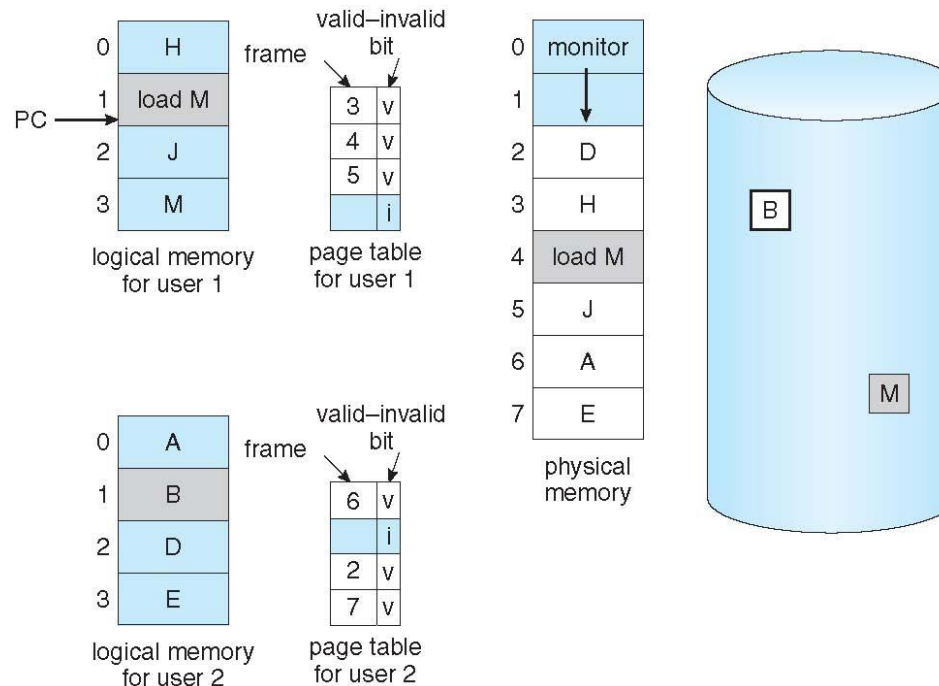- Operation-system examples

# Two major problems

- **Two major problems in demand paging**
  - Page-replacement algorithm
  - Frame-allocation algorithms

- **Even slight improvement can yield large gain in performance.**
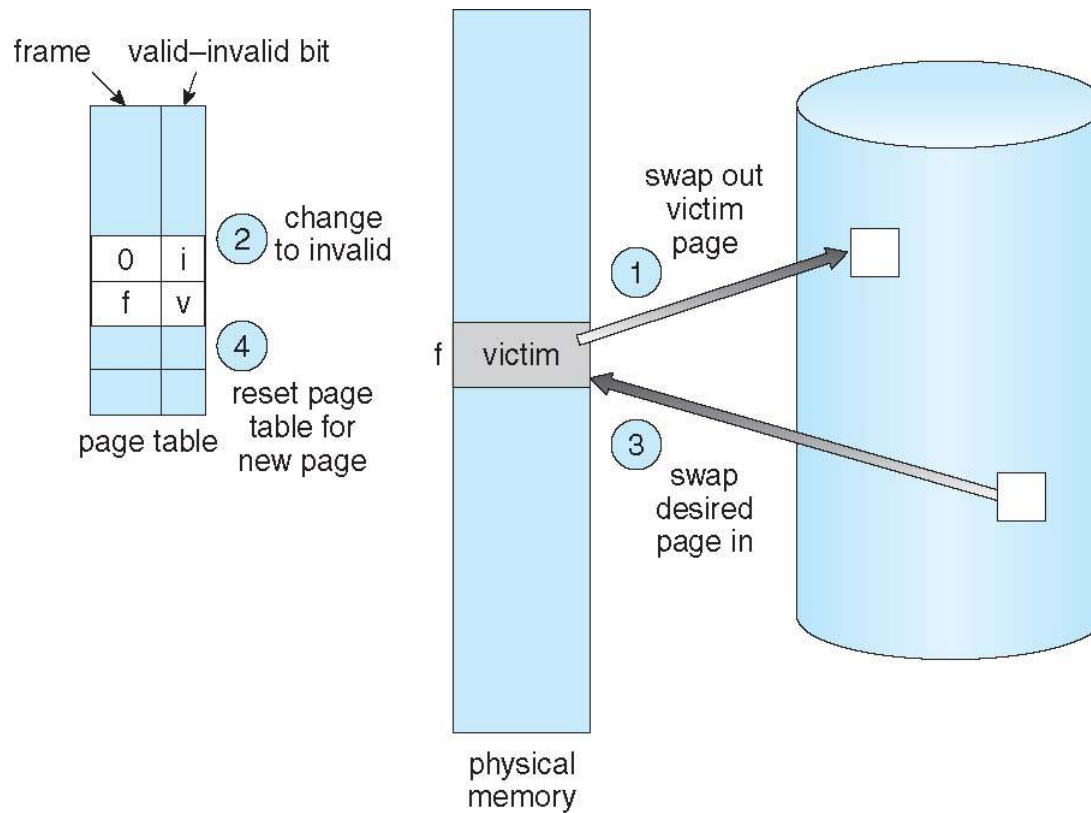  - Effective access time = $(1-p)*ma + p * $ <page fault time>

# Page Replacement

■ **Page replacement**
- If no frame is free at a page fault, we find a frame not being used currently, and swap out
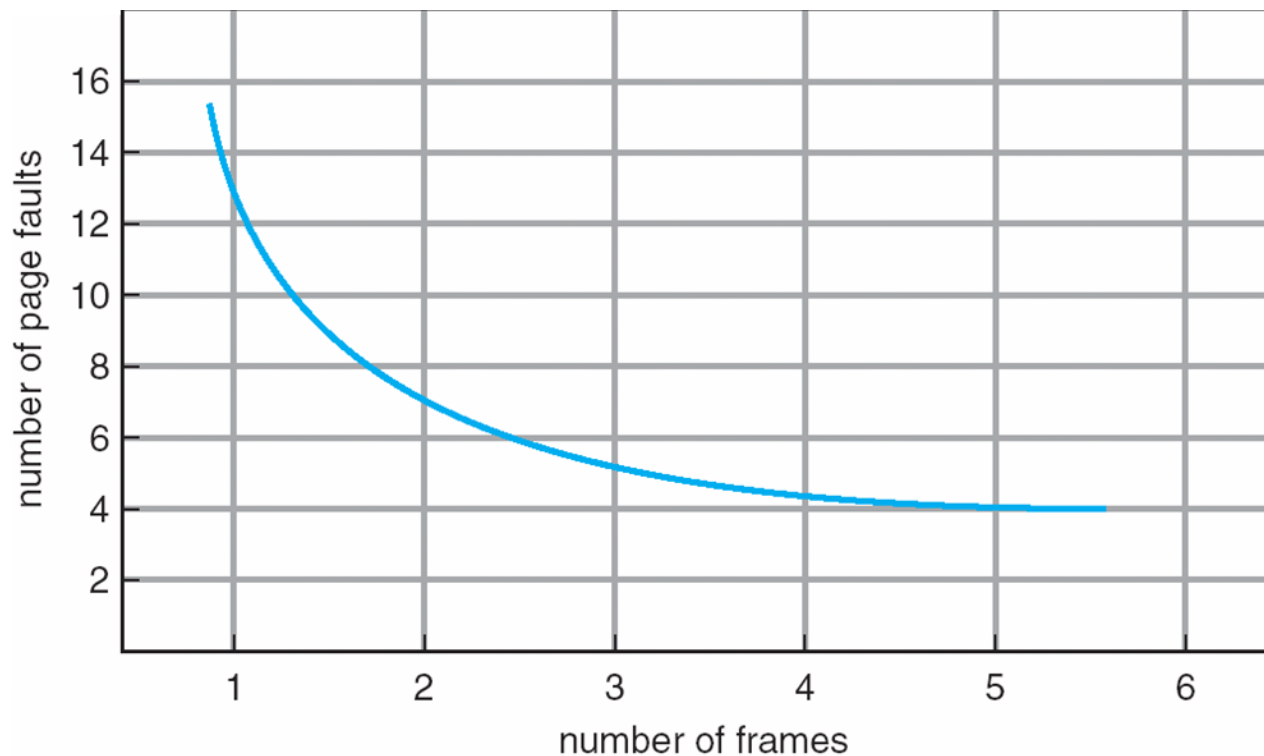- Writing overhead can be reduced by modify-bit (or dirty-bit) for each frame

# Page Replacement

# Page Faults vs. Number of Frames

- ■ In general, the more frames, the fewer page faults

# Page Replacement Algorithms

- FIFO page replacement

- Optimal page replacement (in theory)

- Least-recently-used (LRU) page replacement

- LRU-approximation page replacement

- ETC.

# FIFO Page Replacement

- **First-in, first-out**: when a page should be replaced, the oldest page is chosen.
  - Easy, but not always good

reference string

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

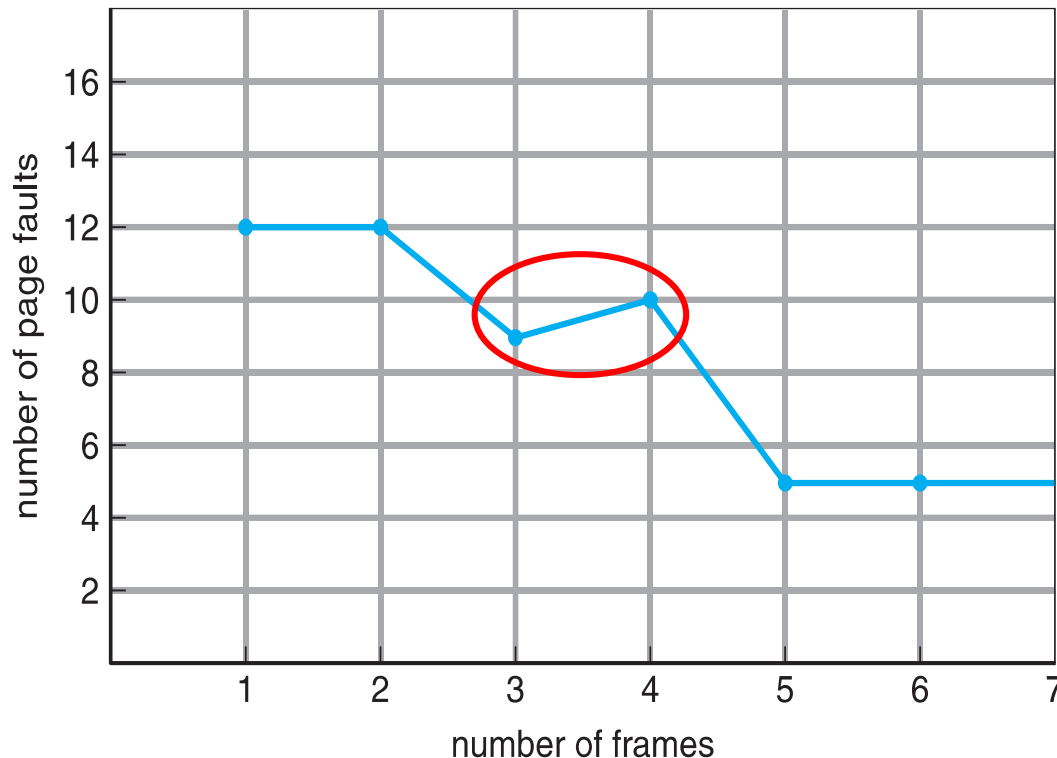| 7 | 7 | 7 | 2 |  | 2 | 2 | 4 | 4 | 4 | 0 |  | 0 | 0 |  | 7 | 7 | 7 |
|   | 0 | 0 | 0 |  | 3 | 3 | 3 | 2 | 2 | 2 |  | 1 | 1 |  | 1 | 0 | 0 |
|   |   | 1 | 1 |  | 1 | 0 | 0 | 0 | 3 | 3 |  | 3 | 2 |  | 2 | 2 | 1 |

page frames

  - # of page faults: 15
- **Problem: Belady's anomaly**

# FIFO Page Replacement

■ **Belady's anomaly**: # of faults for 4 frames is greater than # of faults for 3 frames

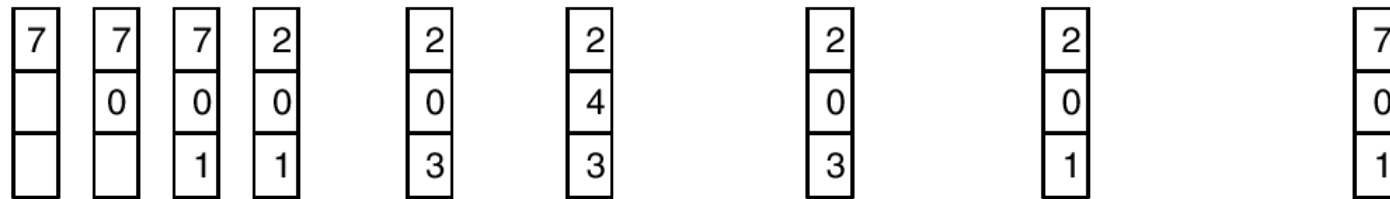(Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5)



Page-fault rate may increase as the number of frames increase.

# Optimal Page Replacement

- **Replace the page that will not be used for the longest period of time**

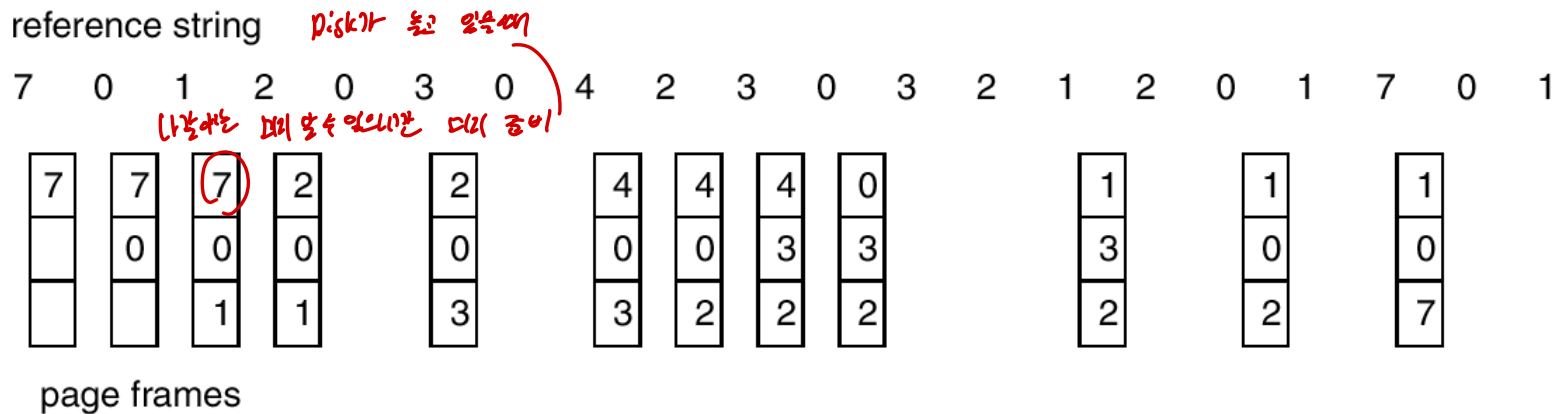reference string

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

page frames

- ▪ # of page faults: 9

- **Problem: It requires future knowledge**

# LRU Page Replacement

■ **LRU (Least Recently Used)**: replace page that has not been used for longest period of time



reference string

Disk가 돌고 있을때

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

나중에 페이지가 있으면 다시 줌어

| 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

- ■ # of page faults: 12
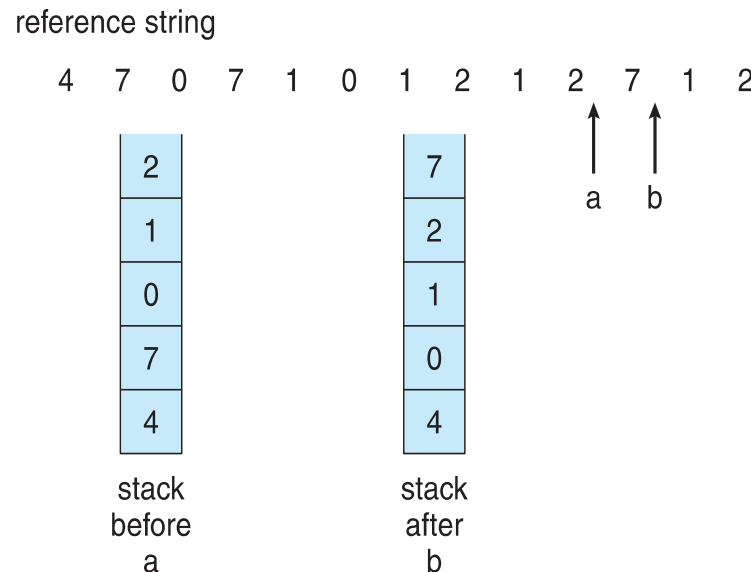- ■ LRU is considered to be good and used frequently

# Implementation of LRU

- **Using counter (logical clock)**
  - Associate with each page-table entry a time-of-used field
  - Whenever a page is referenced, clock register is copied to its time-of-used field

- **Using stack of page numbers**
  - If a page is referenced, remove it and put on the top of the stack

reference string

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

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

stack
before
a

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

stack
after
b

이전도 느림

└ 한번 쓰이면 검색이 필요함

# Stack Algorithm

■ Does LRU cause the Belady's anomaly?

■ Stack algorithm: an algorithm for which the set of pages in memory for $n$ frames is always a subset of the set of pages that would be in memory with $n + 1$ frames.

   ▪ Never exhibit Belady's anomaly
   ▪ LRU is a stack algorithm

$n$ frames

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

$n + 1$ frames

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

# LRU-Approximation Page Replacement

- **Motivation**
  - LRU algorithm is good, but few system provide sufficient supports for LRU
  - However, many systems support reference bit for each page
    - We can determine which pages have been referenced, but not their order.
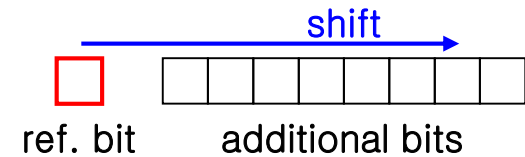
- **LRU-approximation algorithms**
  - Additional-reference-bit algorithm
  - Second-chance algorithm

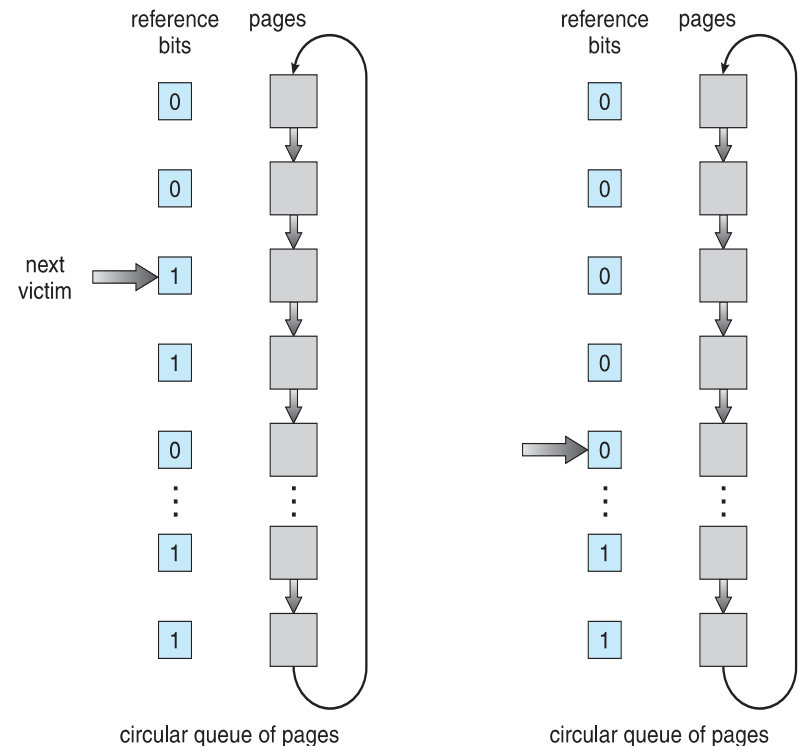# LRU-Approximation Page Replacement

- **Additional-reference-bit algorithm**
  - Gain additional ordering information by recording the reference bits at regular intervals.
    - Ex) keep 8-bit byte for each page and records current state of reference bits of each page

- **Second-chance algorithm**
  - Basically, FIFO algorithm
  - If reference bit of the chosen page is 1, give a second chance
    - The reference bit is cleared.

shift

ref. bit    additional bits



reference bits    pages

next victim → 1

circular queue of pages          circular queue of pages

# Page-Buffering Algorithms

- **Pool of free frames**
  - Some system maintain a list of free frames.
  - When a page fault occurs, a victim frame is chosen from the pool.
  - Frame number of each free page can be kept for next use.

- **Keep a list of free frames and remember which page was in each frame.**
  - The old page can be reused.
  - OS typically applies ZFOD(zero-fill on demand) technique.

- **List of modified pages**
  - Whenever the page device is idle, a modified page is written to disk.

# Allocation of Frames

- **How do we allocate the fixed amount of free memory among various processes?**
  - How many frames does each process get?

- **Minimum number of frames for each process**
  - # of frames for each process decreases
    - -> page-fault rate is increases
    - -> performance degradation
  - Minimum # of frames should be large enough to hold all different pages that any single instruction can reference.

# Allocation Algorithms

- **Equal allocation**
  - Split m frames among n processes
  - -> m/n frames for each process

- **Proportional allocation** 크기에 비례해서 준 말함
  - Allocate available memory to each process according to its size
$$a_i = s_i/S * m$$
    - $a_i$: # of frames allocated to process $p_i$
    - $s_i$: size of process $p_i$
    - $S = \Sigma\ s_i$

  - Variation: frame allocation based on ···
    - Priority of process
    - Combination of size and priority

# Global vs. Local Allocation

process의 경계 생각 X , 더 효율적인 선택 가능, 남의 공간도 자기 쓰는

- **Global replacement**: a process can select a replacement frame from the set of all frames, including frames allocated to other processes
    - A process cannot control its own page-fault rate

    내께 3에서 아두가

- **Local replacement**: # of frames for a process does not change
    - Less used pages of memory can't be used by other process

-> global replacement is more common method.

# Agenda

- Background
- Demand paging
- Copy-on-Write (COW)
- Page replacement
- Allocation of frames
- <u>Thrashing</u>
- Memory-mapped files
- Allocating kernel memory
- Other considerations
- Operation-system examples

# Thrashing

- If a process does not have enough frames to support pages in active use, it quickly faults again and again.

- A process is thrashing if it is spending more time paging than executing.

# Thrashing → process 별로 Thrashing이 안 일기나 만큼 Memory를 3라야함.

- **If trashing sets in, CPU utilization drops sharply.**
  - Degree of multiprogramming should be decreased

# Thrashing

*memory access가 균일하게 일어나지 않는다.*

*시간에 따라서 메모리의 어떤 영역접근 공간가 많이 쓰인다.*

*access될때 page마다 같이 쓰이는 page가 있다.*

*( 필기 대상는 데이터의 수 > frame 수*
*⤷ page fault rate이 급격하게 증가*

*최근의 memory access 공간*

- To prevent thrashing, a process must be provided with as many frames as it needs.
  - -> How to know how many frames it needs?

- Locality model
  - Locality: set of pages actively used together
  - A program is generally composed of several localities

# Working-Set Model

- Working set: set of pages in the most recent Δ page references
  - Parameter Δ: working-set window

page reference table

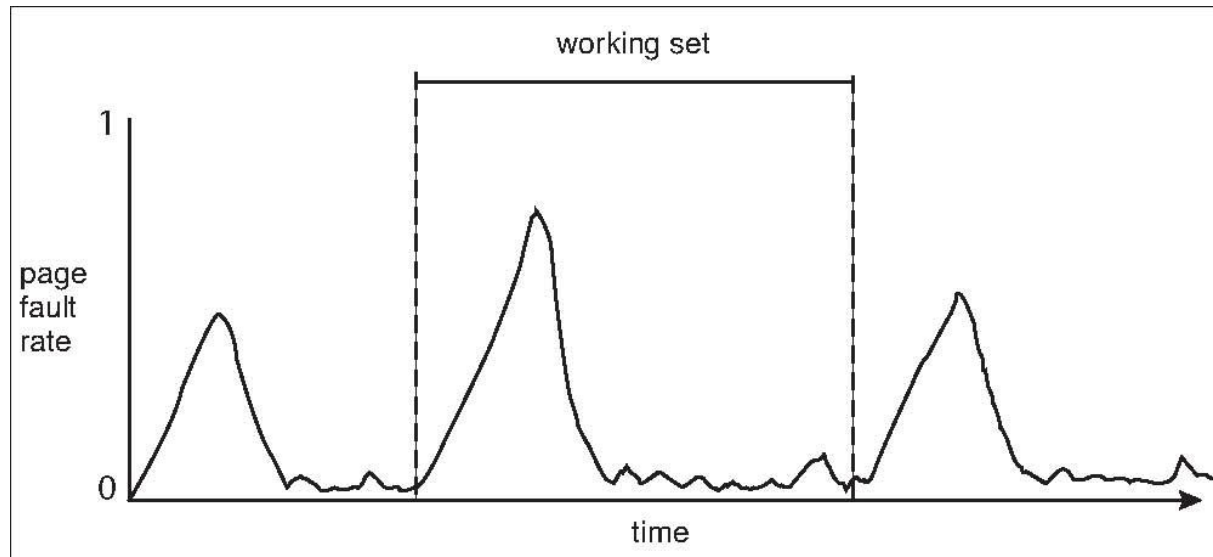. . . 2 6 1 5 7 7 7 7 5 1 6 2 3 4 1 2 3 4 4 4 3 4 3 4 4 4 1 3 2 3 4 4 4 3 4 4 4 . . .

$$\Delta \qquad\qquad\qquad \Delta$$

$t_1 \qquad\qquad\qquad t_2$

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

  - $WSS_i$: working set size of process $p_i$
  - Process $p_i$ needs $WSS_i$ frames
- If total demand is greater than # of available frames, thrashing will occur.

# Working Sets and Page Fault Rates

- Direct relationship between working set of a process and its page-fault rate
- Working set changes over time
- Peaks and valleys over time

# Page−Fault Frequency

■ **Alternative method to control trashing: control degree of multiprogramming by page−fault frequency (PFF)**

　　■ If PFF of a process is too high, allocate more frame

　　■ If PFF of a process is too low, remove a frame from it

# Agenda

- Background
- Demand paging
- Copy-on-Write (COW)
- Page replacement
- Allocation of frames
- Thrashing
- <u>Memory-mapped files</u>
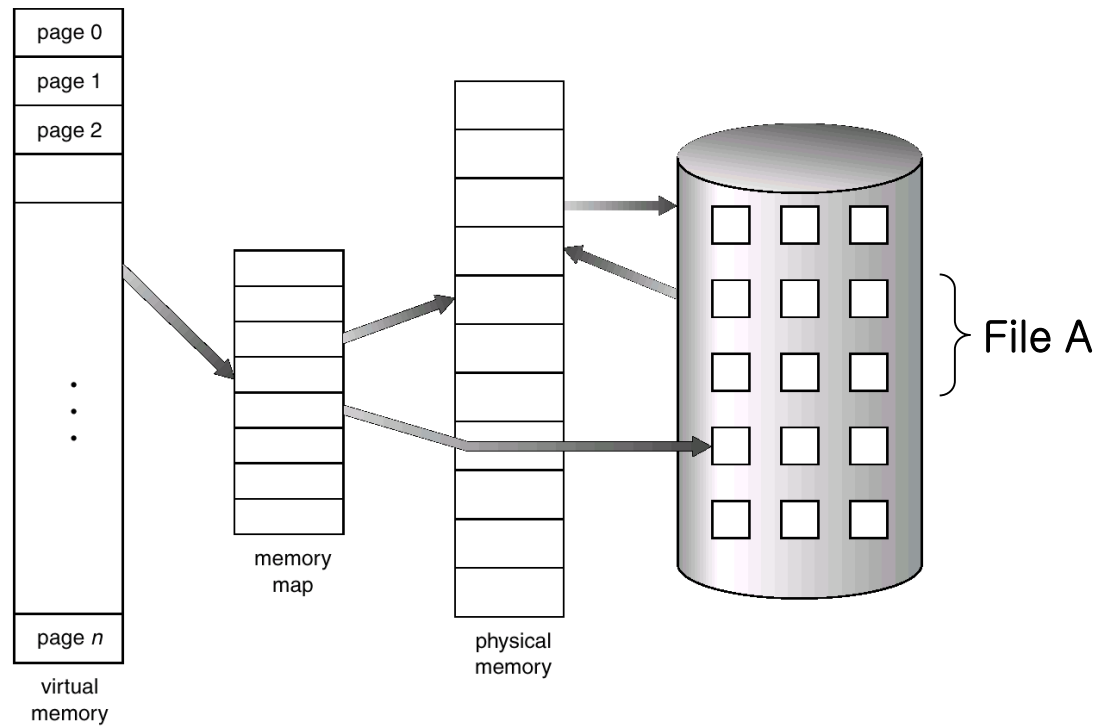- Allocating kernel memory
- Other considerations
- Operation-system examples
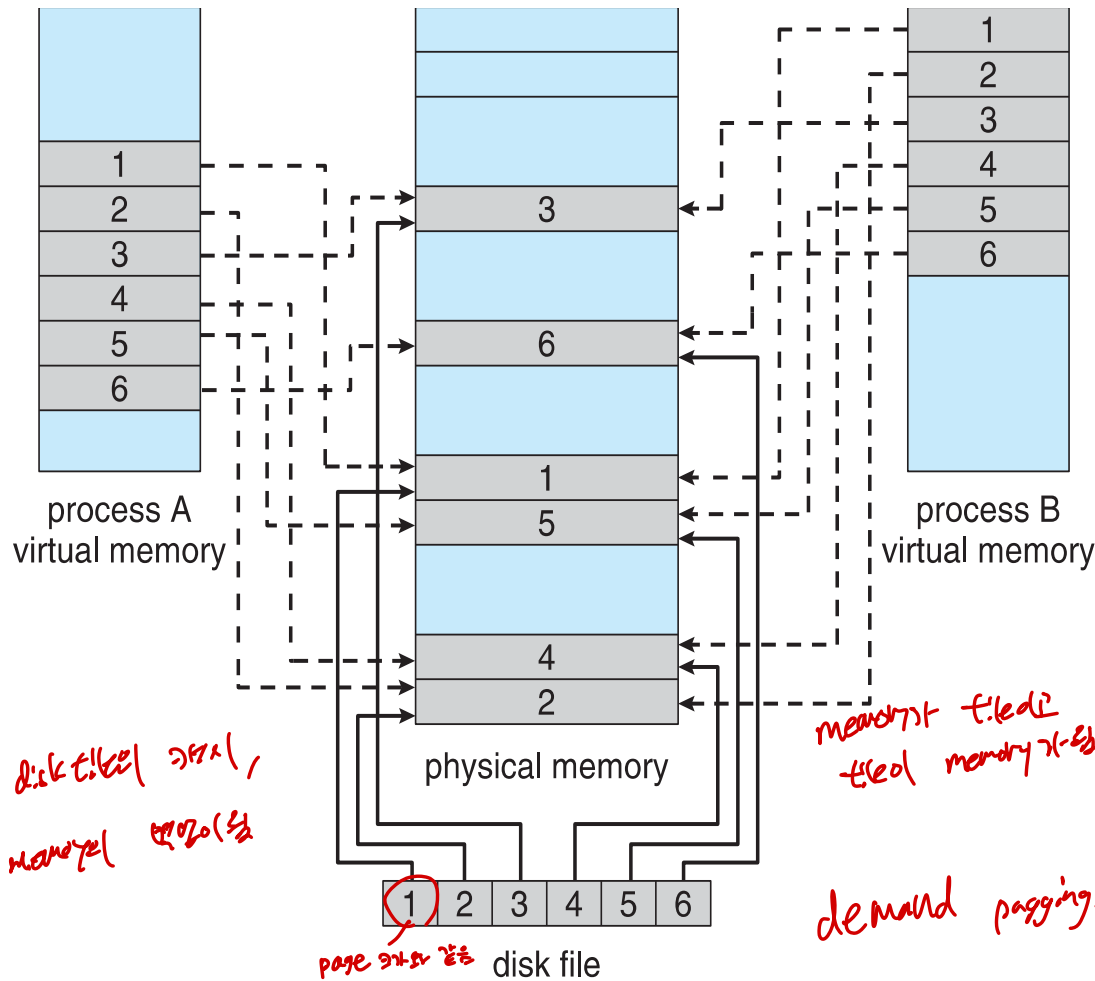
# Memory-Mapped Files

- **Memory-mapped file**: a part of virtual address space is logically associated with a file
  - Map a disk block to a page in memory
  1. Initial access proceeds through demand paging results in page fault.
  2. A page-sized portion of the file is read from the file system into a physical page.
  3. Subsequent access is through memory access routine.
  4. When the file is closed, memory-mapped data are written back to disk.

- **Advantages**
  - Reduces overhead of read() and write()
  - File sharing

# Memory-Mapped Files

# Memory–Mapped Files



Memory–mapped file
Shared by two processes

# Shared Memory in Win32 API

- **Sharing memory-mapped file is similar to shared memory**
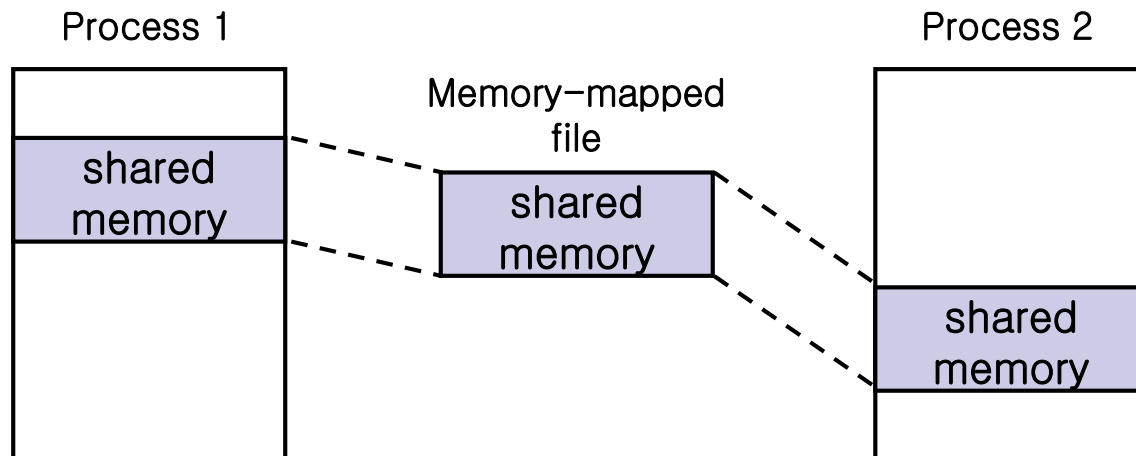    - On POSIX and Windows, shared memory is accomplished by memory mapping files

# Shared Memory in Win32 API

■ **Process 1**

1. Create a file to be shared
   - CreateFile(…)

2. Create file mapping (named shared memory object)
   - CreateFileMapping(…)

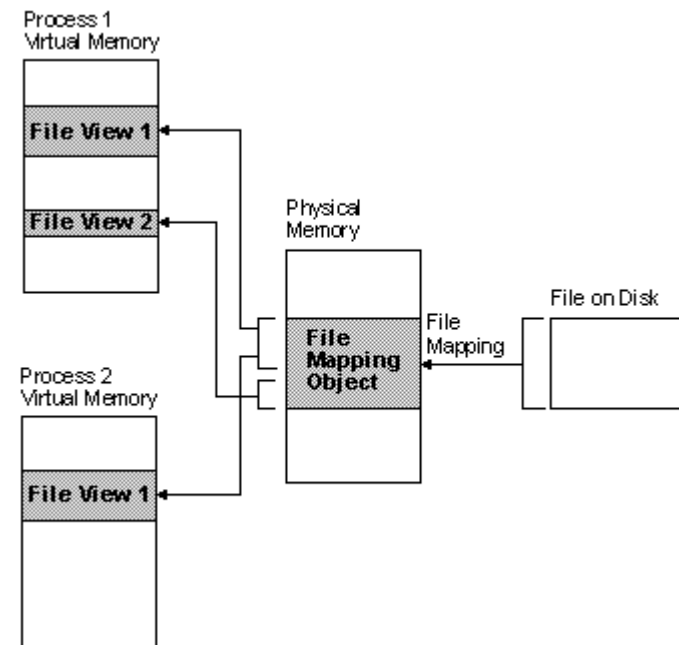3. Establish a view of mapped file in virtual address space
   - MapViewOfFile(…)

■ **Process 2**

1. Open file mapping
   - OpenFileMapping(…)

2. Establish a view of mapped file in virtual address space
   - MapViewOfFile(…)

Process 1

Memory-mapped file

Process 2

shared memory

shared memory

shared memory

# Memory Mapped File on Windows

■ **File mapping is the association of a file's contents with a portion of the virtual address space of a process.**

  ■ The system creates a file mapping object to maintain this association.

  ■ A file view is the portion of virtual address space that a process uses to access the file's contents.

  ■ It also allows the process to work efficiently with a large data file.

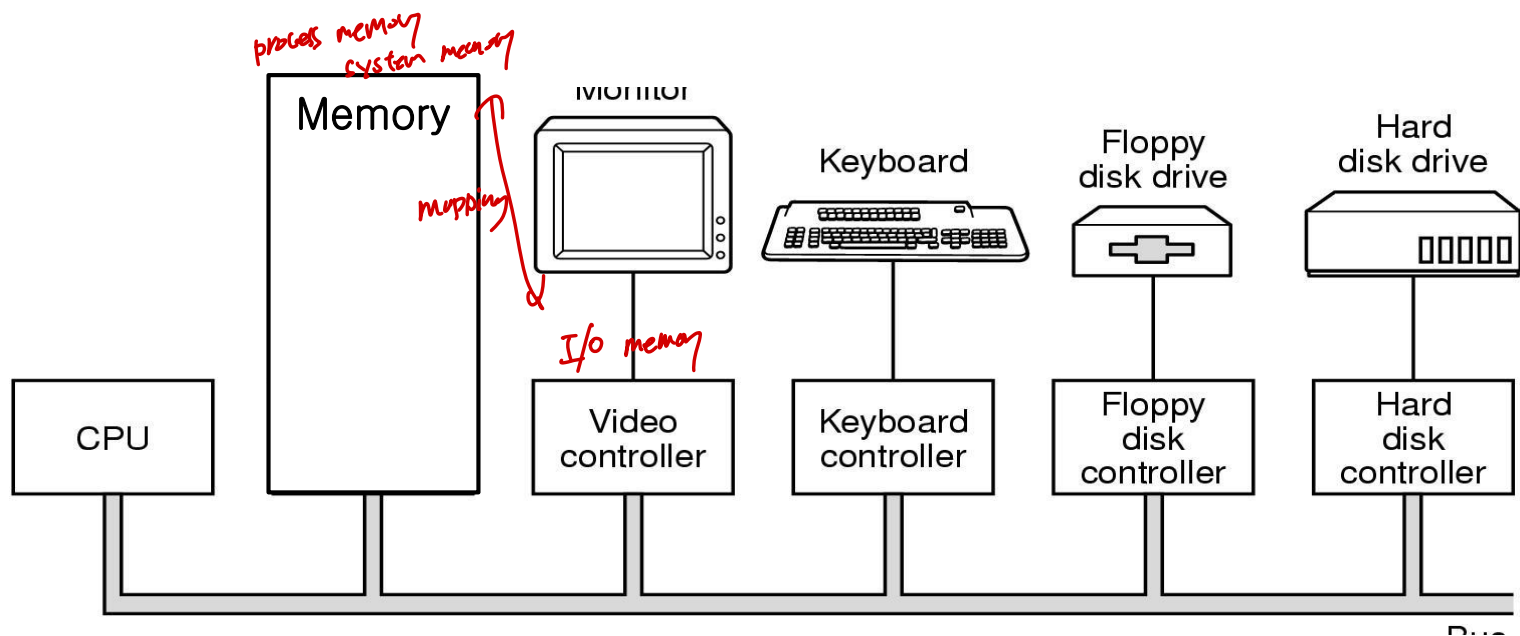  ■ Multiple processes can also use memory-mapped files to share data.

# Reading Assignment

- **Visit and study the following web pages to learn memory mapped file on Windows**
  - File mapping
    - https://msdn.microsoft.com/en-us/library/windows/desktop/aa366556(v=vs.85).aspx

  - File mapping functions
    - https://msdn.microsoft.com/en-us/library/windows/desktop/aa366781(v=vs.85).aspx#file_mapping_functions

# Memory-Mapped I/O

- **I/O devices are accessed through ⋯**
  - Device registers in I/O controller to hold command and data
  - Usually, special instructions transfer data between device registers and memory

# Memory-Mapped I/O

■ **Memory-mapped I/O**: ranges of memory addresses are set aside and mapped to the device registers

- In IBM PC, each location on screen is mapped to a memory location
- Serial/parallel ports – data transfer through reading/writing device registers (ports)

■ **Programmed I/O (PIO)**

Ex) Sending a long string of bytes through memory-mapped I/O

- CPU sets a byte to a register and set a bit in the control register to signal that the data is available.
- Device receives the data and clears the bit in the control register to signal CPU that it is ready for the next byte.

cf. Interrupt driven I/O

# Agenda

- Background
- Demand paging
- Copy-on-Write (COW)
- Page replacement
- Allocation of frames
- Thrashing
- Memory-mapped files
- <u>Allocating kernel memory</u>  *paging 보다 어떻게 효율적*
- <u>Other considerations</u>
- Operation-system examples
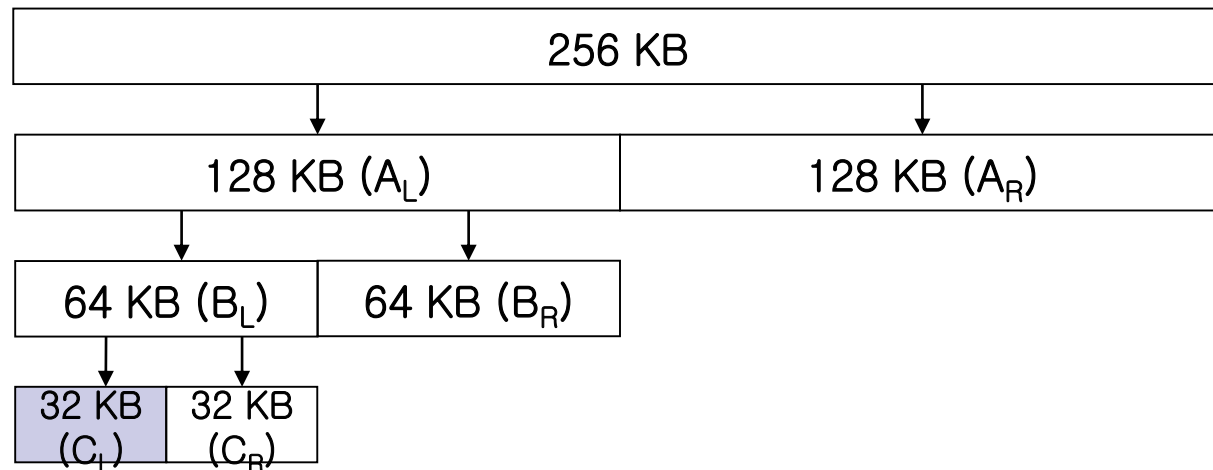
# Allocating Kernel Memory

- **Allocation of kernel memory requires special handling**
  - Kernel requests memory for data structures of varying sizes
    - Many OS's do not subject kernel code/data to the paging system
  - Certain H/W devices interact directly with physical memory
    - Memory should reside in physically contiguous pages.

- **Strategies for kernel memory allocation**
  - Buddy system
  - Slab allocation

# Buddy System

- **Buddy system**: allocates memory from a fixed-size segment consisting of physically contiguous pages
  - Power-of-2 allocator

  Ex) initially 256 KB is available, 21 KB was requested

| 256 KB |
|--------|

| 128 KB ($A_L$) | 128 KB ($A_R$) |
|----------------|----------------|

| 64 KB ($B_L$) | 64 KB ($B_R$) |
|---------------|---------------|

| 32 KB ($C_L$) | 32 KB ($C_R$) |
|---------------|---------------|

  - Advantage: <u>easy to combine adjacent buddies</u>  *frac 최소 표현하여*
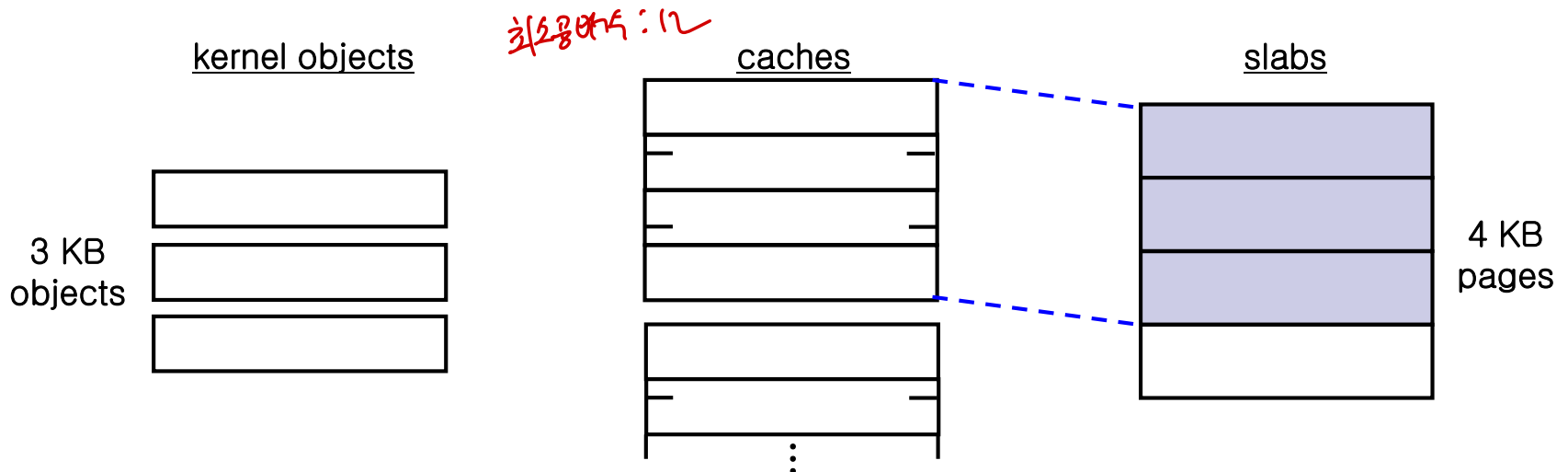  - Disadvantage: internal fragmentation  $12 < 32$

# Slab Allocation

*struct 를 단위 → 크기 일정*

*ex) PCB, TCB, shared memory block*

- **Motivation: mismatch between allocation size and requested size**
  - Page-size granularity vs. byte-size granularity
  - Applied since Solaris 2.4 and Linux 2.2

- **Cache for each unique kernel data structure**    *똑같은 크기의   structure가 할당/해제됨*
  - A slab is made up of one or more physically contiguous pages
  - A cache consists of one or more slabs

kernel objects     *최소할당수: 12*     caches     slabs



3 KB objects
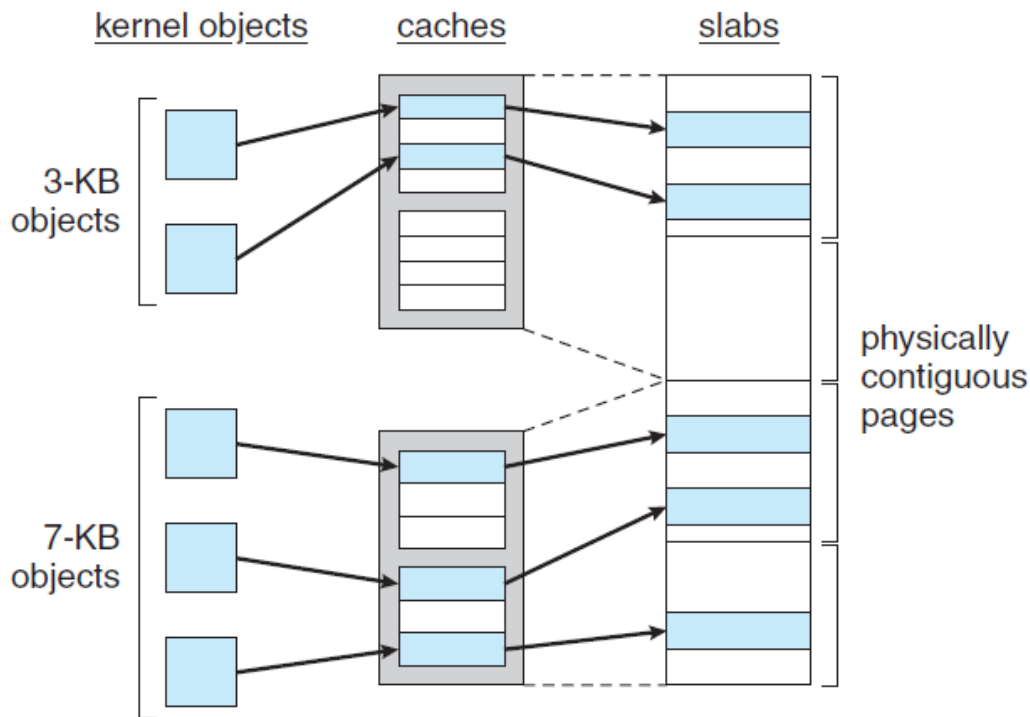
4 KB pages

# Slab Allocation

- **Single cache is for each unique kernel data structure**
  - Each cache filled with objects – instantiations of the data structure.

    Ex) cache for process descriptor, cache for file objects, cache for semaphore, ⋯

- **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 is allocated.

# Slab Allocation

- **Benefits**
  - No memory waste due to fragmentation
  - Memory requests can be satisfied quickly
  - ➔ Suitable for data structures that are allocated and deallocated frequently.

# Other Considerations

- Prepaging
- Page size
- TLB reach
- Inverted page tables
- Program structure
- I/O interlock

# Prepaging  *page faults를 줄여보고*

*메모(예측해서)*

- ## A problem of pure demand paging: a large number of page faults

- ## Prepaging: bring all pages that will be needed at one time to reduce page faults.

   Ex) working-set model *(같이 다니는 page들의 set)*

   - Important issue: cost of prepaging vs. cost of servicing corresponding page faults

# Page Size 증가하는 추세

intel page size $\begin{cases} 4K \\ 4M \end{cases}$

- ## Issues about page size

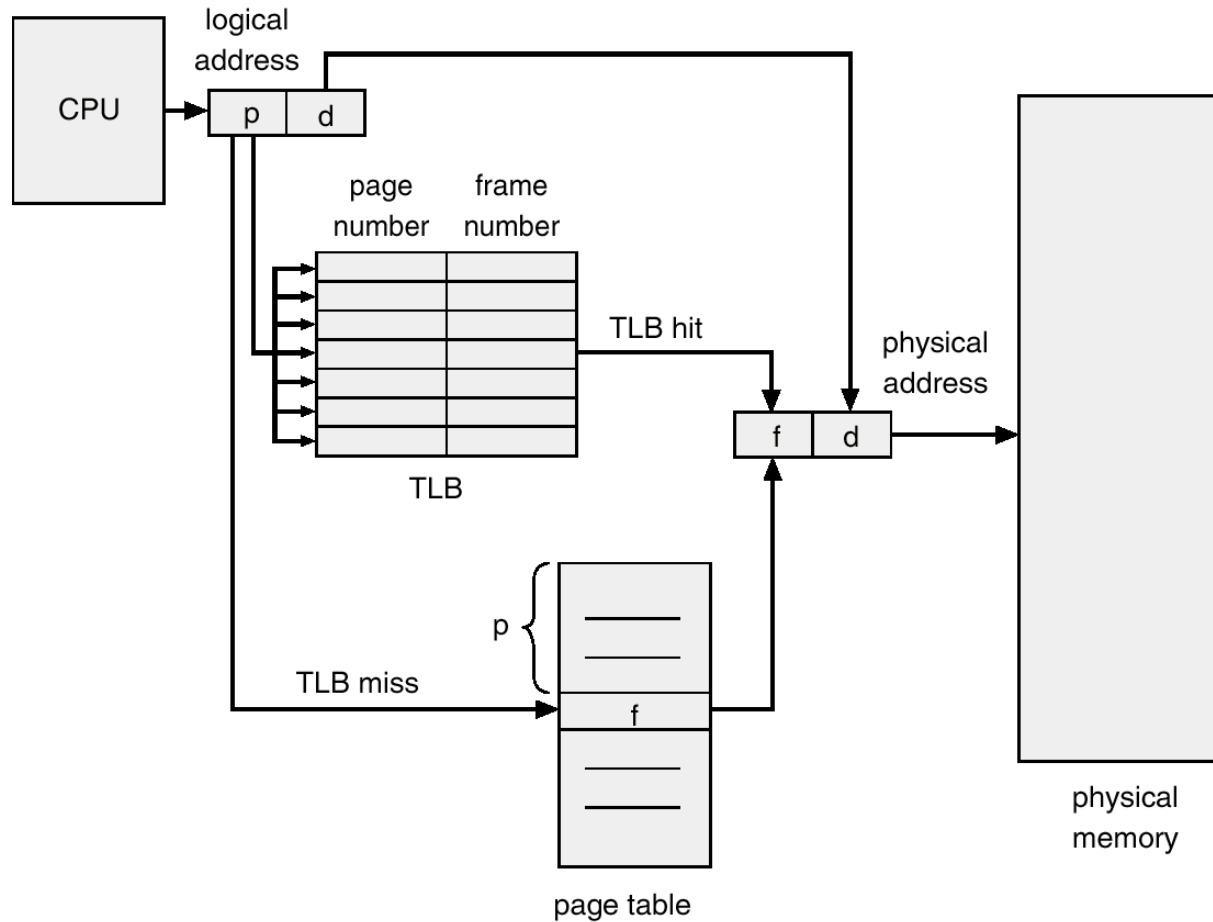|  | smaller page | larger page |
|---|---|---|
| Size of page table | large | small |
| Memory utilization | better | worse |
| I/O latency | large | small |
| Locality | good | bad |
| Page fault | many | few |

Memory 으로
↳ overhead

page 찾기
13시 대한

한번 나올때 cost가큼

- Historical trend: page size is getting larger

# TLB Reach

- To improve TLB hit ratio, size of TLB should be increased.

-> but associate memory is expensive, power hungry

- TLB reach: amount of memory accessible from TLB
  - TLB reach = <# of entries in TLB> * <page size>

- TLB reach can increase by increasing page size
  - However, with large page, fragmentation also increases.
  
  ➔ S/W managed TLB (OS support several different page sizes)
  
  Ex) UltraSparc, MIPS, Alpha
  
  Cf) PowerPC, Pentium: H/W managed TLB

# TLB Reach

# Program Structure

- User don't' have to know about nature of memory. But, if user knows underlying demand paging, performance can be improved

  Ex) If page size is 128 words, B is better than A

  ```
  for(j = 0; j < 128; j++)
      for(i = 0; i < 128; i++)          A
          data[i][j] = 0;
  ```
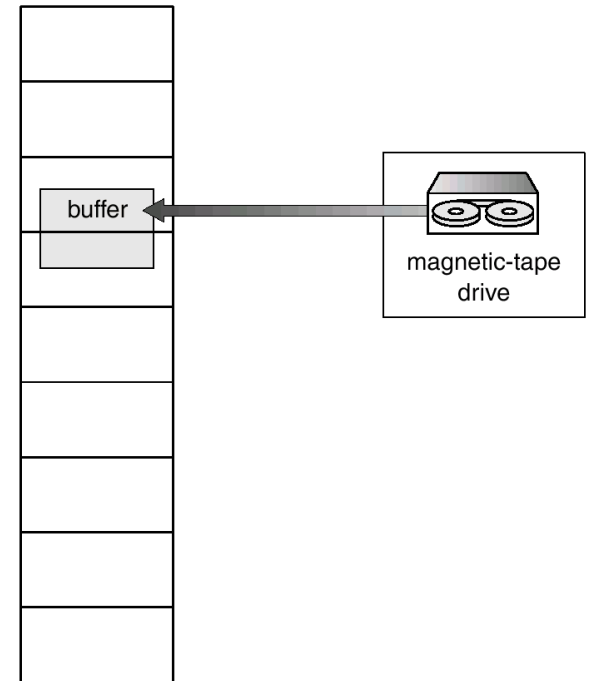
  ```
  for(i = 0; i < 128; i++)
      for(j = 0; j < 128; j++)          B
          data[i][j] = 0;
  ```

# I/O Interlock

- **Memory space related to I/O transfer should not be replaced**

- **Remedies**
  - I/O transfer is performed only through system memory
  - Locking pages in memory

buffer

magnetic-tape drive

# Agenda

- Background
- Demand paging
- Copy-on-Write (COW)
- Page replacement
- Allocation of frames
- Thrashing
- Memory-mapped files
- Allocating kernel memory
- Other considerations
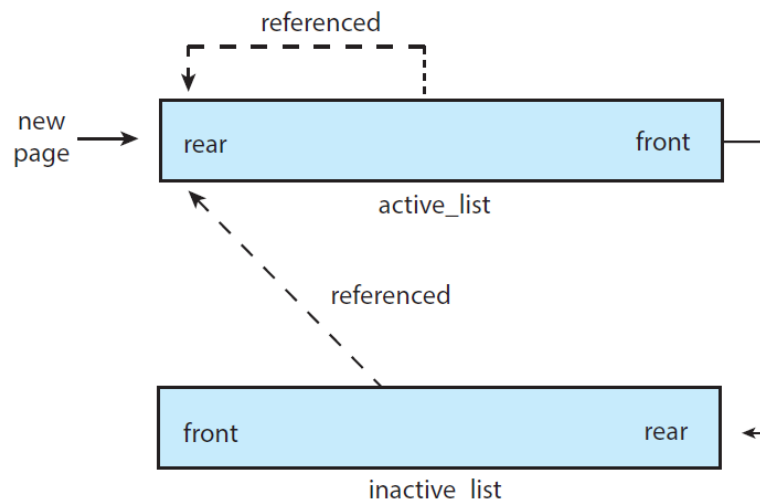- <u>Operation-system examples</u>

# Linux

- **Demand paging**
  - Allocates pages from a list of free frames.
  - A global page-replacement policy similar to the LRU-approximation clock algorithm

- **Accessed bit**
  - Each page has an accessed bit that is set whenever the page is referenced.
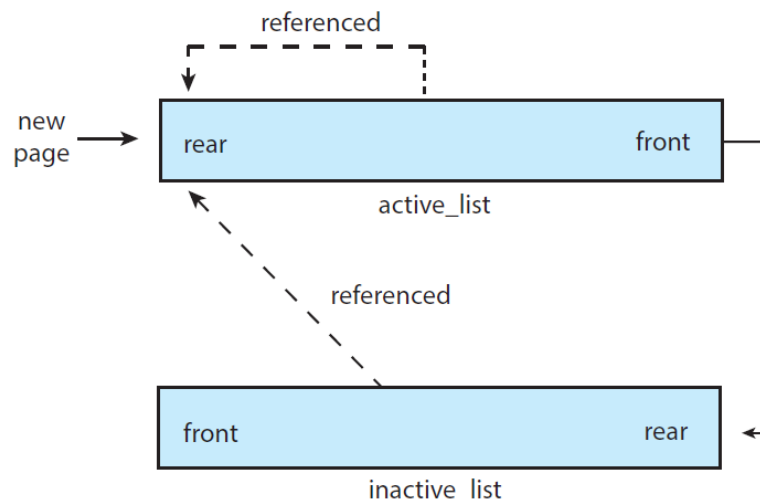  - Periodically, the accessed bits for pages in the active list are reset

# Linux

- **Linux maintains two types of page lists**
  - Active list contains the pages that are considered in use
    - Over time, the least recently used page will be at the front of the active list.
  - Inactive list contains pages that have not recently been referenced and are eligible to be reclaimed.
    - If a page in the inactive list is referenced, it moves back to the rear of the active list

# Linux

- **Kernel swap daemon process kswapd**
  - A page-out daemon process
  - Periodically awakens and checks the amount of free memory in the system
  - If free memory falls below a certain threshold, kswapd begins scanning pages in the inactive list and reclaiming them for the free list

# Windows

- ## Demand paging with clustering
  - At page fault, Windows brings not only the fault page but also several following pages

- ## Working-set minimum/maximum
  - Working-set minimum (default=50): the virtual memory manager attempts to keep at least this much memory resident in the process whenever the process is active

  - Working-set maximum (default: 345): the virtual memory manager attempts to keep no more than this much memory resident in the process whenever the process is active and available memory is low.