

Chapter 9 Virtual Memory

Contents

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

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

Background

- Normal method for memory management
 - One time
 - Stay in memory for ever
- Actually, the entire program is not needed all the time:
 - The code to handle unusual conditions
 - Arrays, lists and tables
 - Certain options and features of a program

Background

□ Problems

- Big program
- Occupy memory

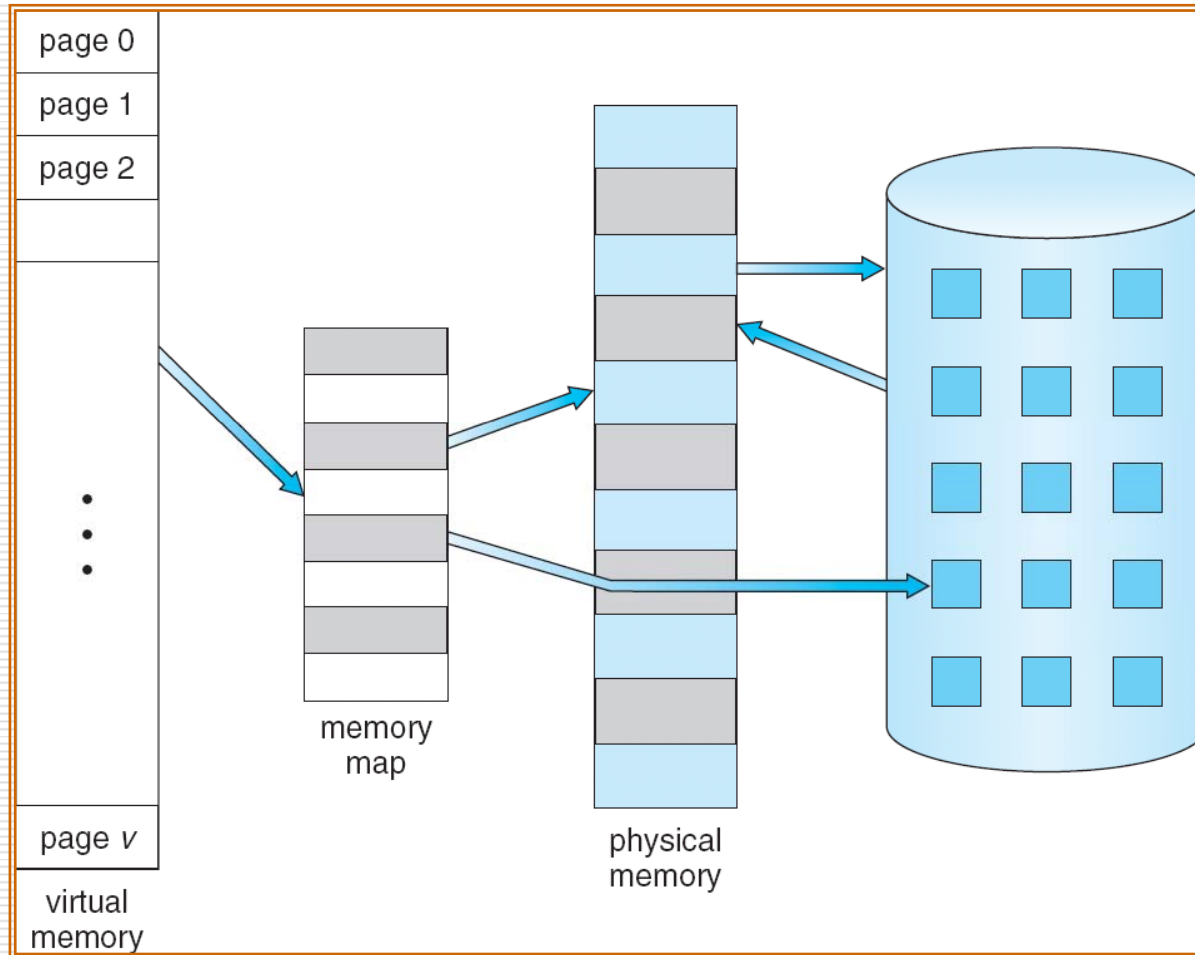
□ Solution

- Increase your memory capacity
- Other technologies: dynamic load, overlap, swap
- Virtual memory

Background

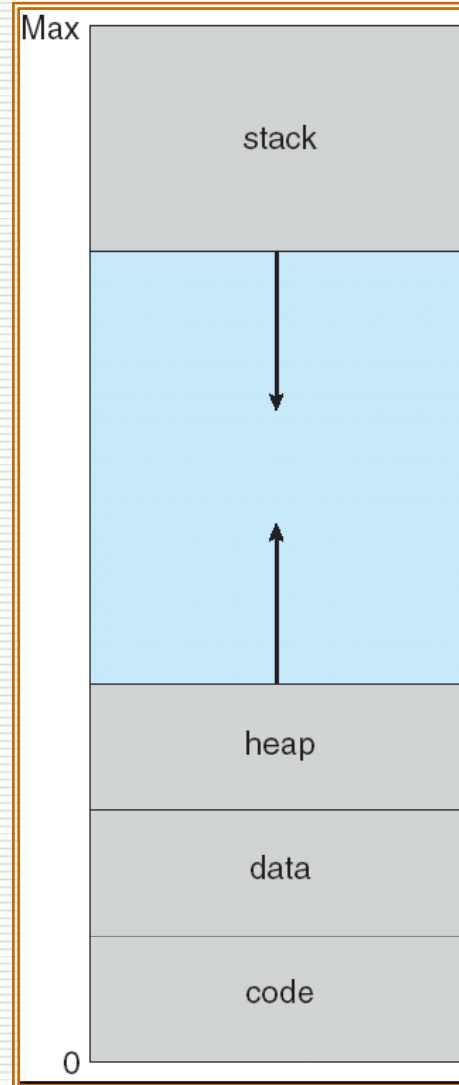
- **Virtual memory** – separation of user logical memory from physical memory.
 - Only part of the program needs to be in memory for execution
 - Logical address space can therefore be much larger than physical address space
 - Allows address spaces to be shared by several processes
 - Allows for more efficient process creation

Virtual Memory That is Larger Than Physical Memory

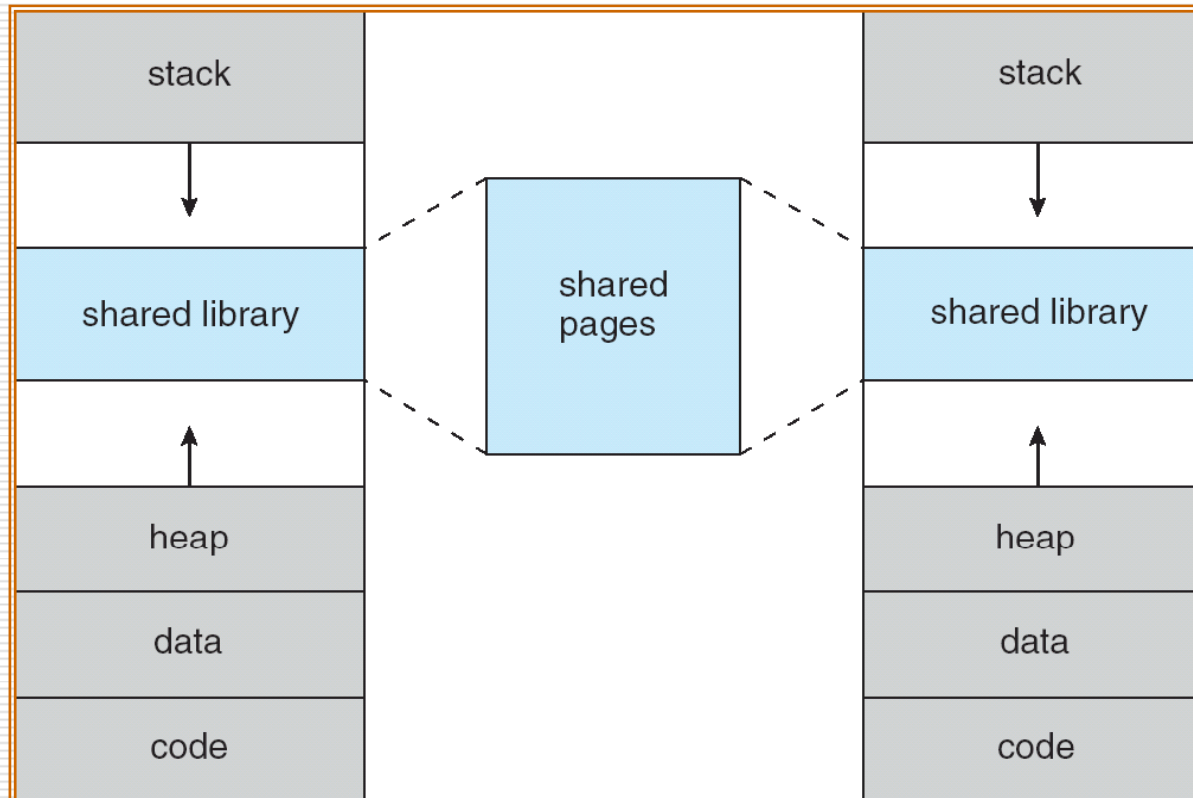


Virtual-address Space

- ❑ It refers to the logical view of how a process is stored in memory
- ❑ It begins at a certain logical address and exists in contiguous memory.



Shared Library Using Virtual Memory



Technologies to implement VM

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

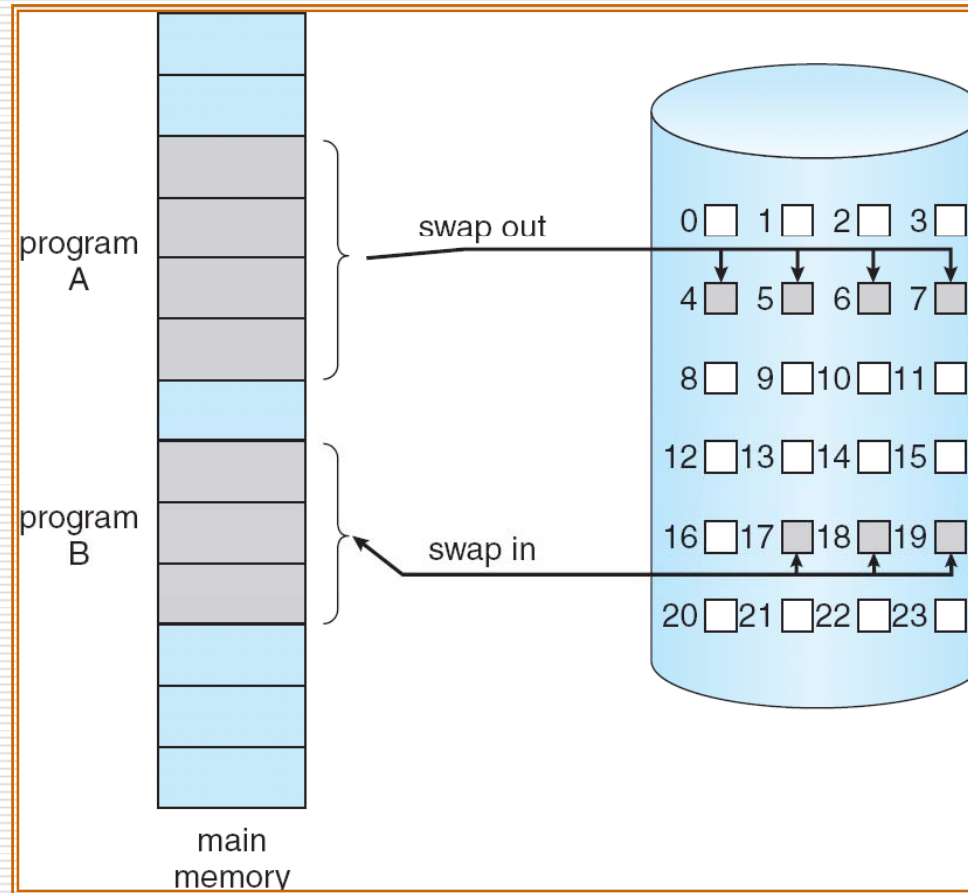
Demand Paging

- Bring a page into memory only when it is needed
 - Less I/O needed
 - Less memory needed
 - Faster response
 - More users

- 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**

Transfer of a Paged Memory to Contiguous Disk Space



Valid-Invalid Bit

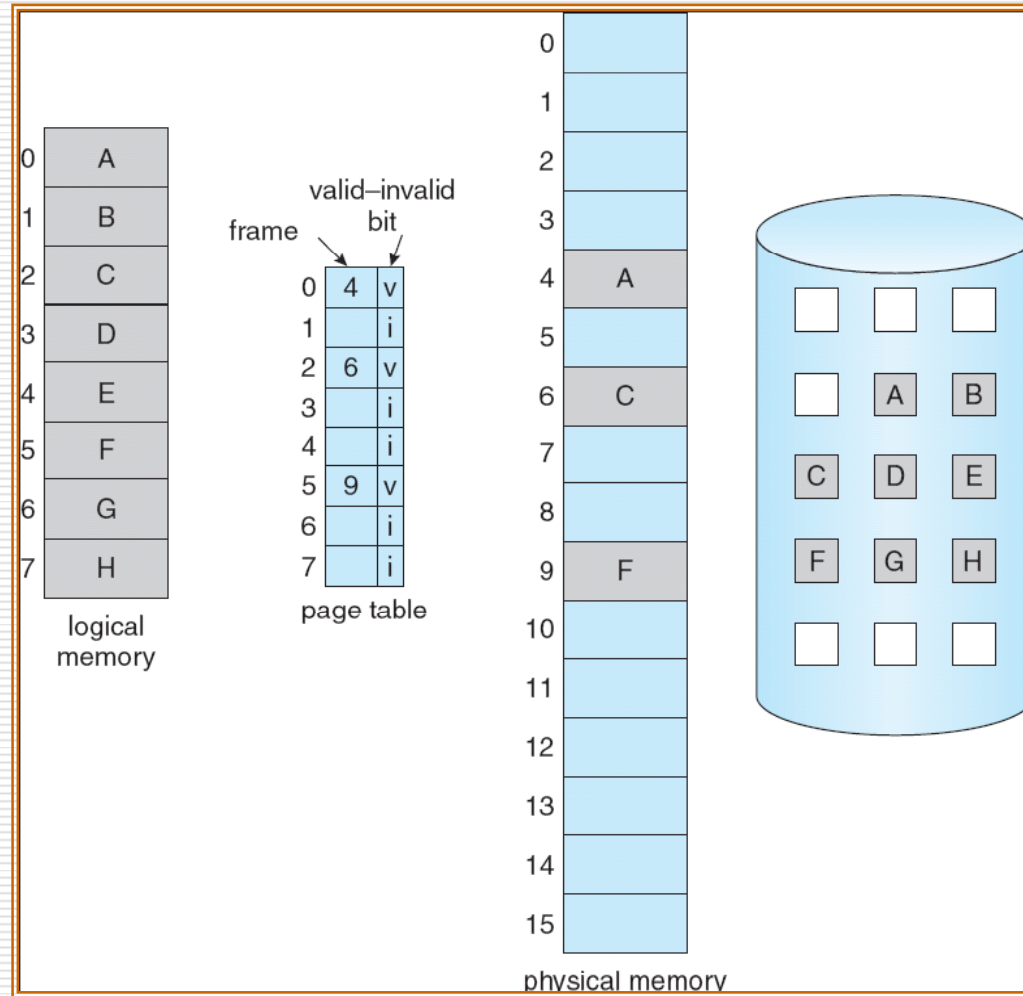
- With each page table entry a valid–invalid bit is associated (**v** \Rightarrow in-memory, **i** \Rightarrow not-in-memory)
- Initially valid–invalid bit is set to **i** on all entries
- Example of a page table snapshot:

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

page table

- During 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



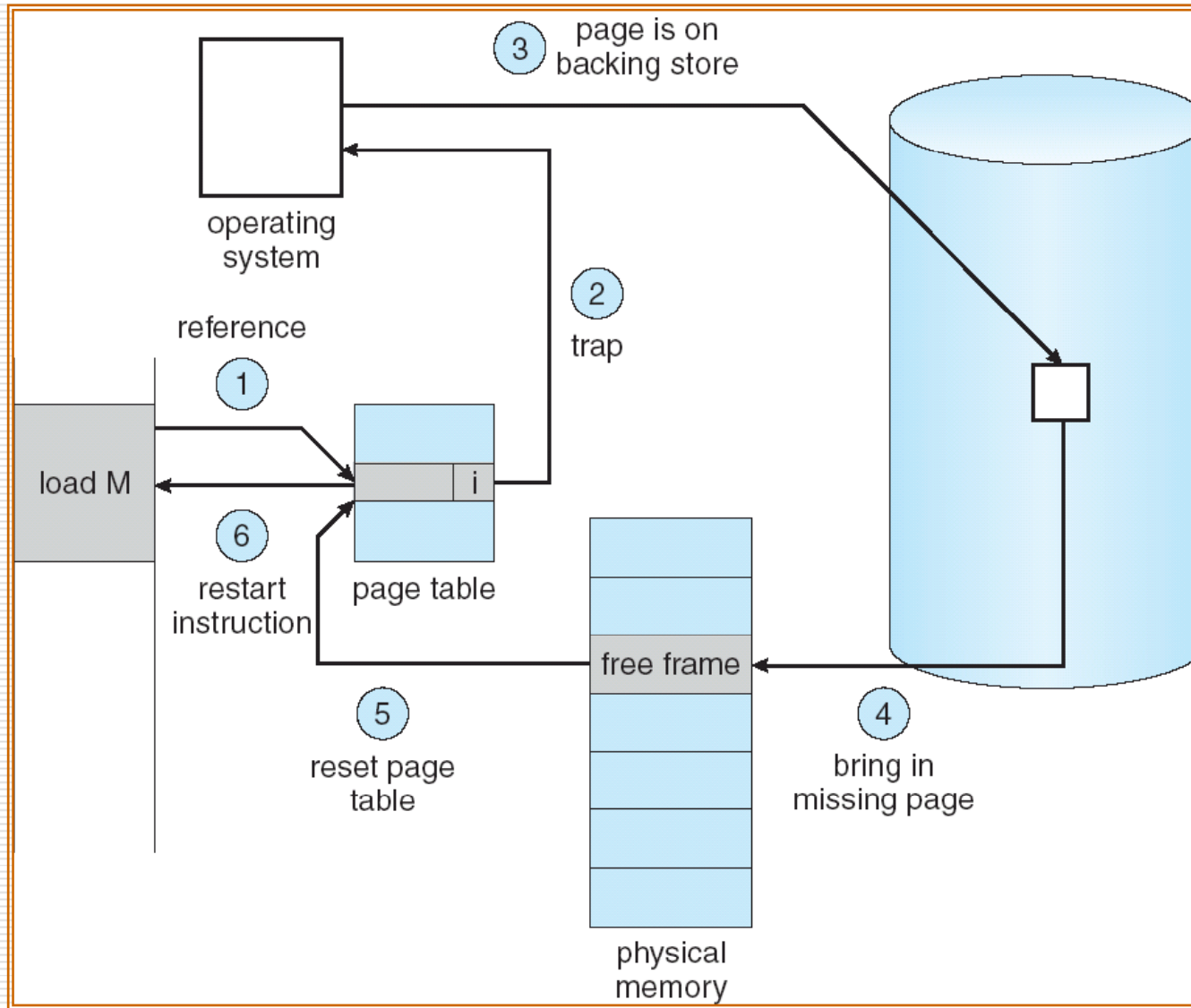
Page Fault

- 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 \Rightarrow abort
 - Just not in memory
2. Get empty frame
3. Swap page into frame
4. Reset tables
5. Set validation bit = **v**
6. Restart the instruction that caused the page fault

Steps in Handling a Page Fault



Demand paging

- Pure demand paging
 - Never bring a page into memory until it is required.
- Locality of reference
- The hardware to support demand paging
 - Page table
 - Secondary memory

Performance of Demand Paging

- Page Fault Rate $0 \leq p \leq 1.0$
 - 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
 - + restart overhead
 -)

Demand Paging Example

- Memory access time = 200 nanoseconds
- Average page-fault service time = 8 milliseconds
- $EAT = (1 - p) \times 200 + p (8 \text{ milliseconds})$
 $= (1 - p) \times 200 + p \times 8,000,000$
 $= 200 + p \times 7,999,800$
- If one access out of 1,000 causes a page fault, then
EAT = 8.2 microseconds.

Process Creation

- Virtual memory allows other benefits during process creation:
 - Copy-on-Write
 - Memory-Mapped Files (later)

Copy-on-Write

- ❑ Copy-on-Write (COW) allows both parent and child processes to initially *share* the same pages in memory

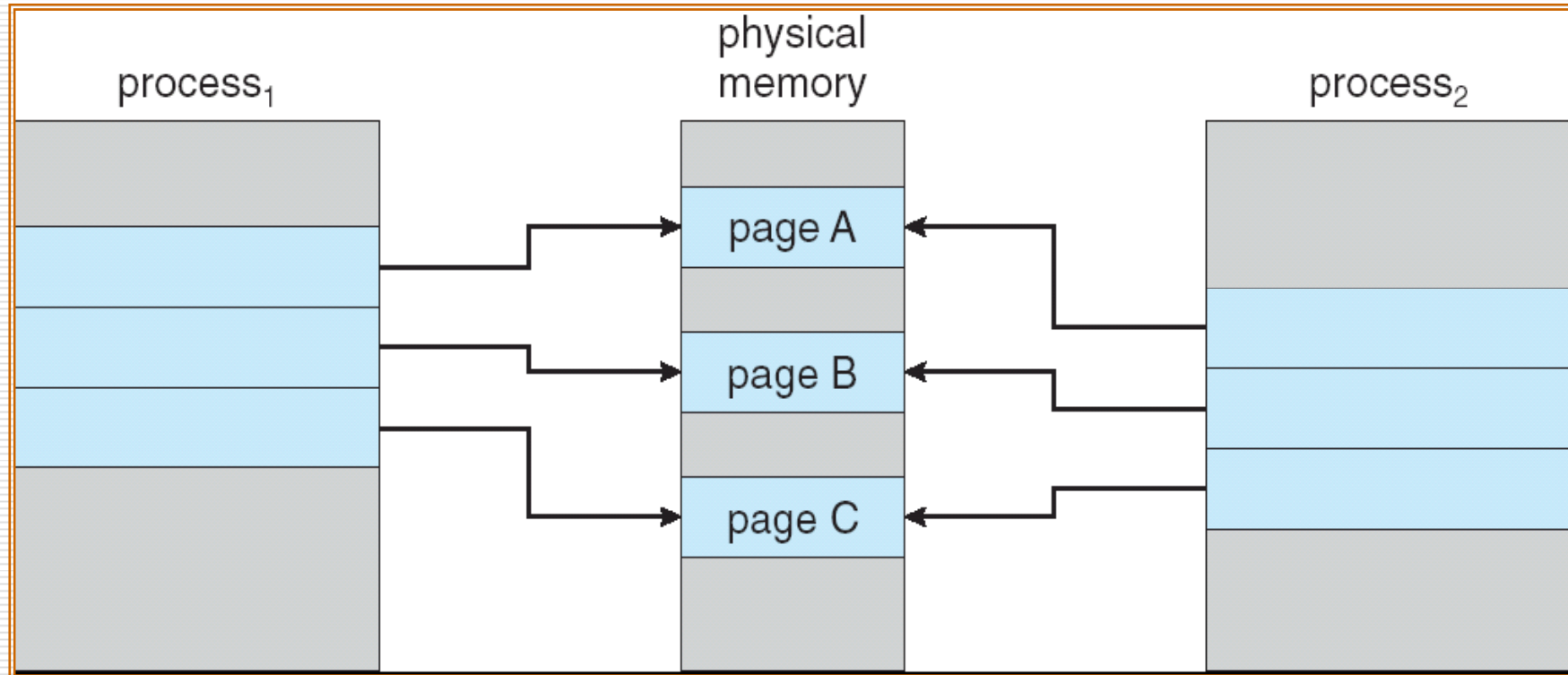
If either process modifies a shared page, only then is the page copied

- ❑ COW allows more efficient process creation as only modified pages are copied
- ❑ Free pages are allocated from a **pool** of zeroed-out pages

vfork

- UNIX系统的某些版本提供fork调用的变种，比如：vfork。
- Vfork产生的子进程会使用父进程的所有地址空间

Before Process 1 Modifies Page C



After Process 1 Modifies Page C

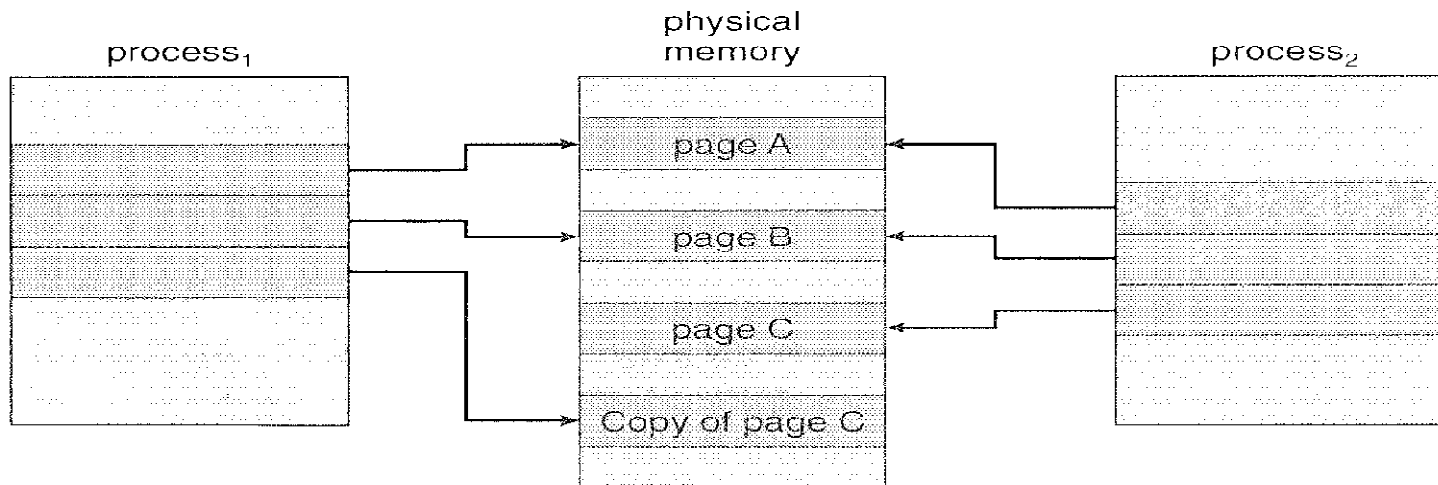
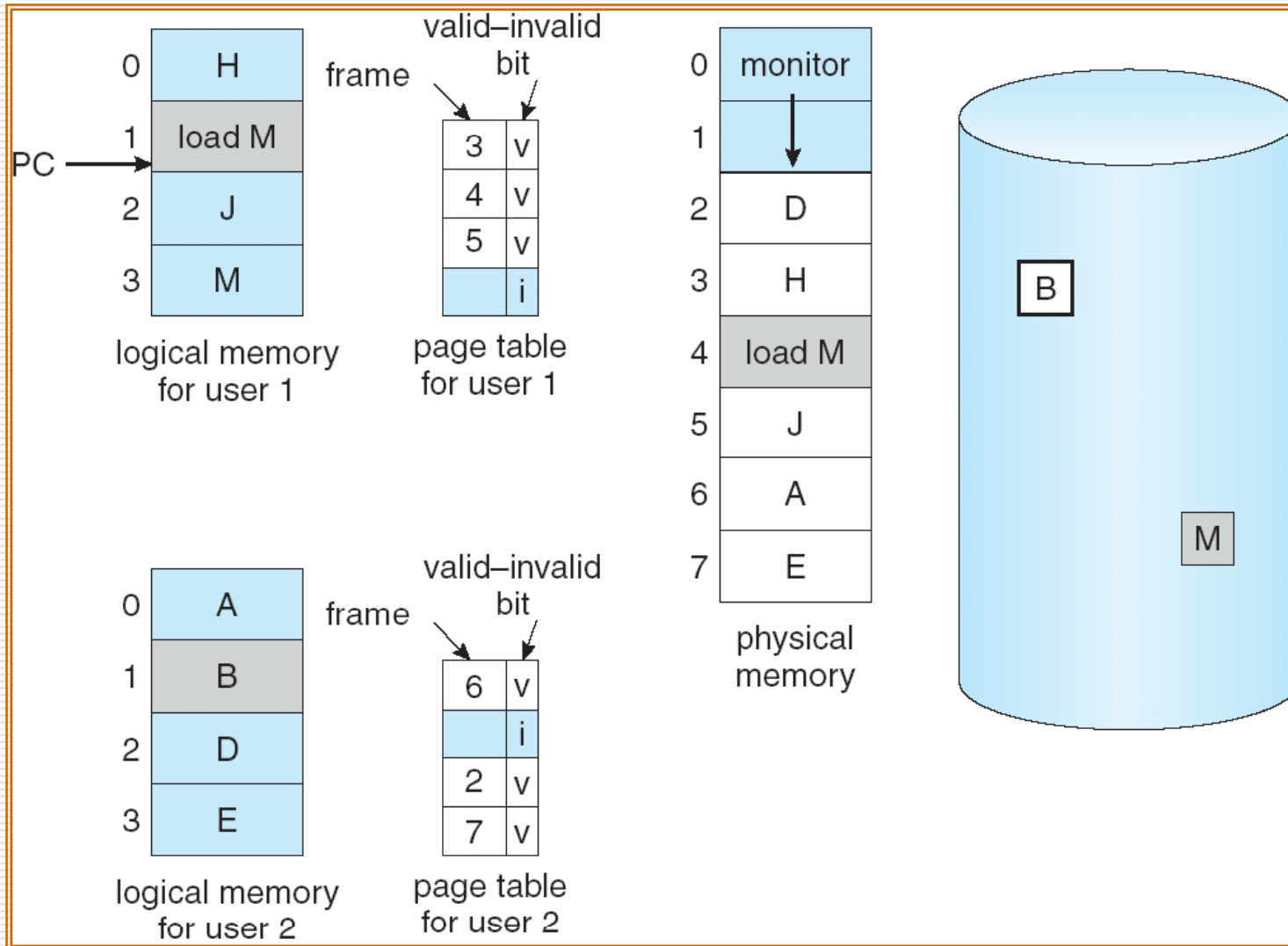


Figure 9.8 After process 1 modifies page C.

What happens if there is no free frame?

- Page replacement – find some page in memory, but not really in use, swap it out
 - algorithm
 - performance – want an algorithm which will result in minimum number of page faults
- Same page may be brought into memory several times

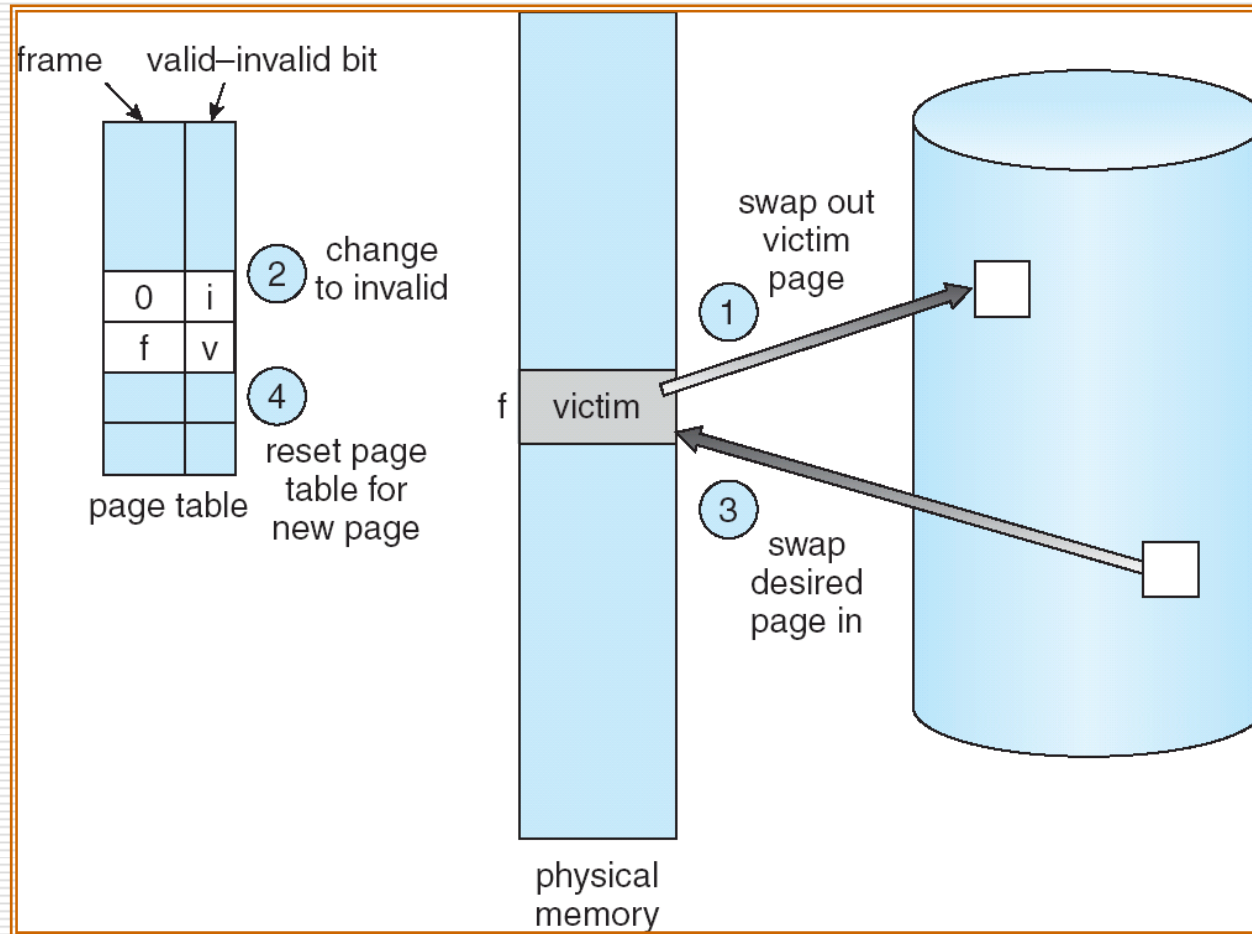
Need For Page Replacement



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
3. Bring the desired page into the (newly) free frame; update the page and frame tables
4. Restart the process

Page Replacement



Page Replacement

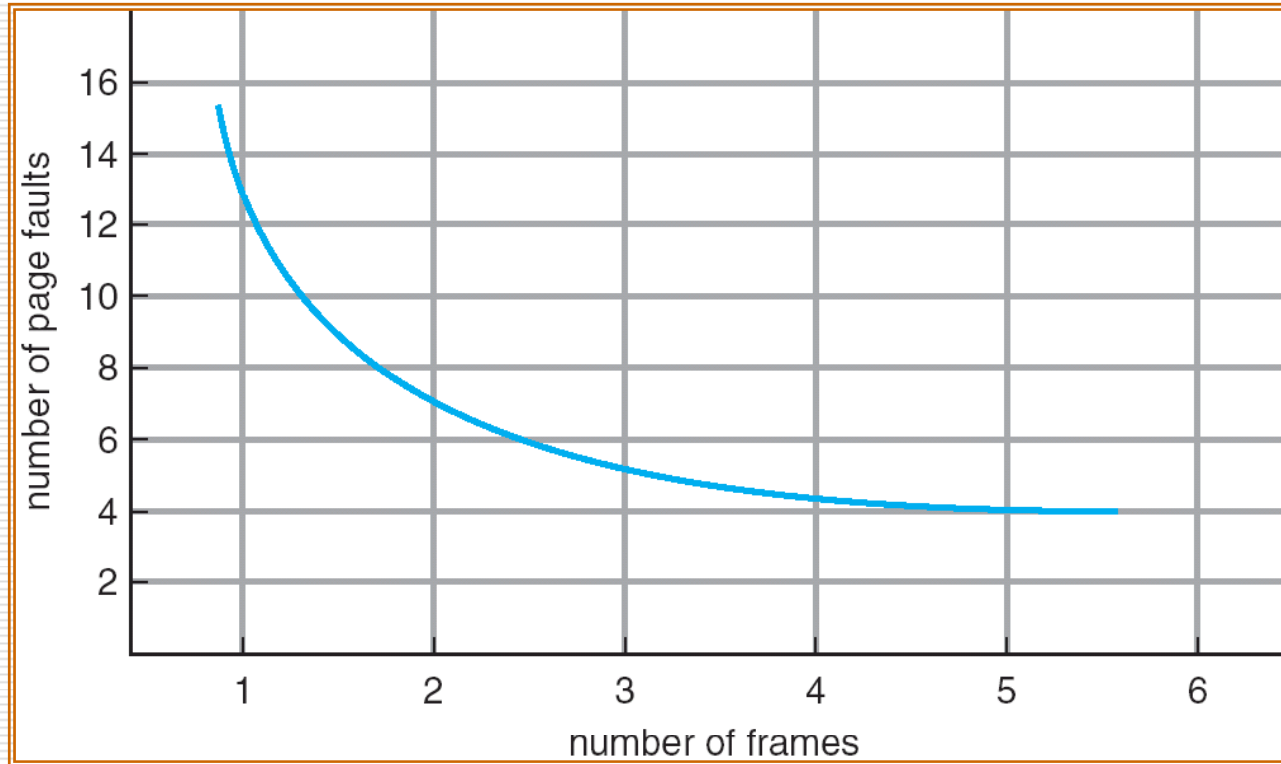
- ❑ Prevent over-allocation of memory by modifying page-fault service routine to include page replacement
- ❑ Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk
- ❑ Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory
- ❑ Two major problems to implement demand paging
 - Frame-allocation algorithm
 - Page-replacement algorithm

Page Replacement Algorithms

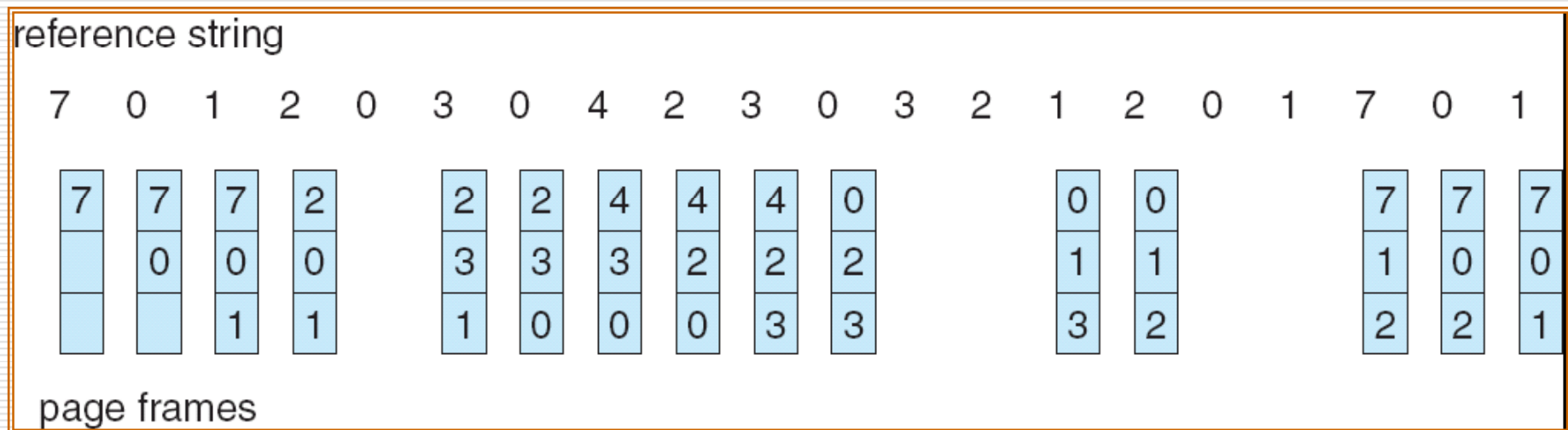
- Want lowest page-fault rate
- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
- In all our examples, the reference strings are as follows:

1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5
7,0,1,2,0,3,0,4,2,3,0,3,2,1,2,0,1,7,0,1

Graph of Page Faults Versus The Number of Frames



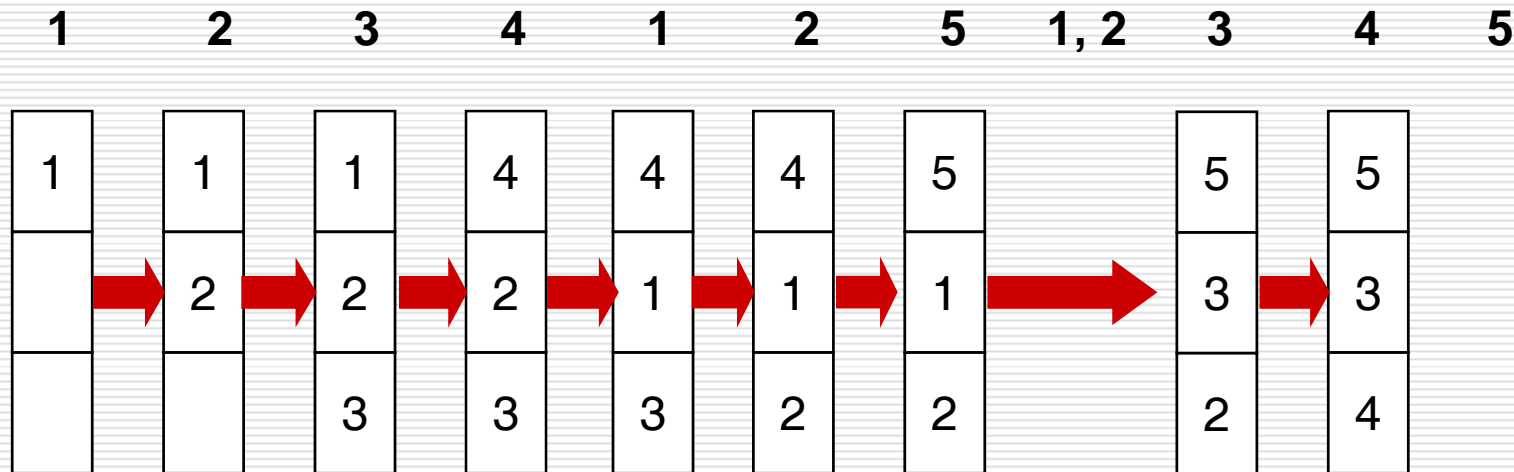
FIFO Page Replacement



15 page faults

FIFO Algorithm-- Belady's Anomaly

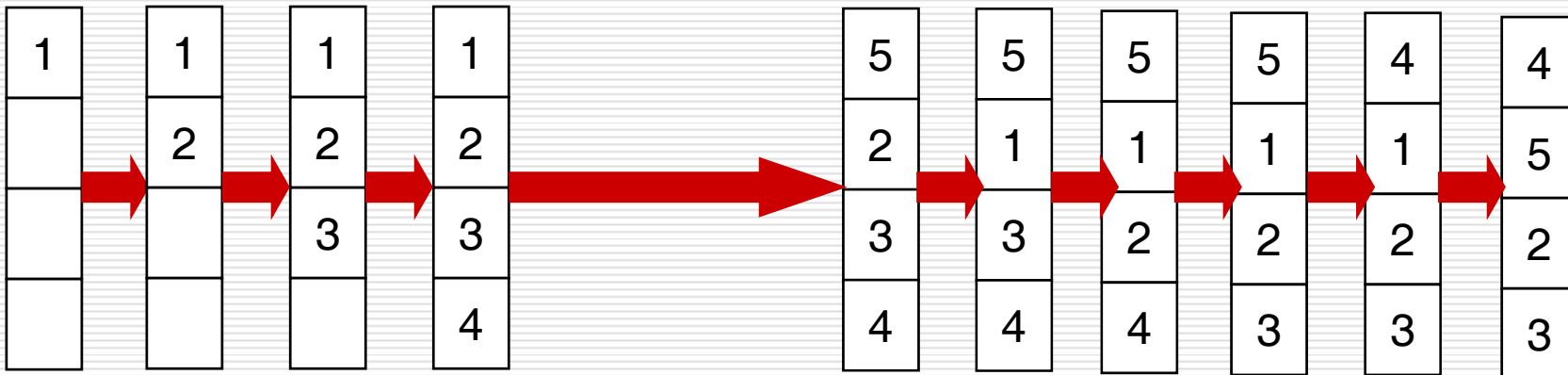
- 在FIFO算法中，有时候帧数的增加反而会使缺页次数增加，如下例: number of frames is 3 or 4
- Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5
- 3 frames (3 pages can be in memory at a time per process)



9 page faults

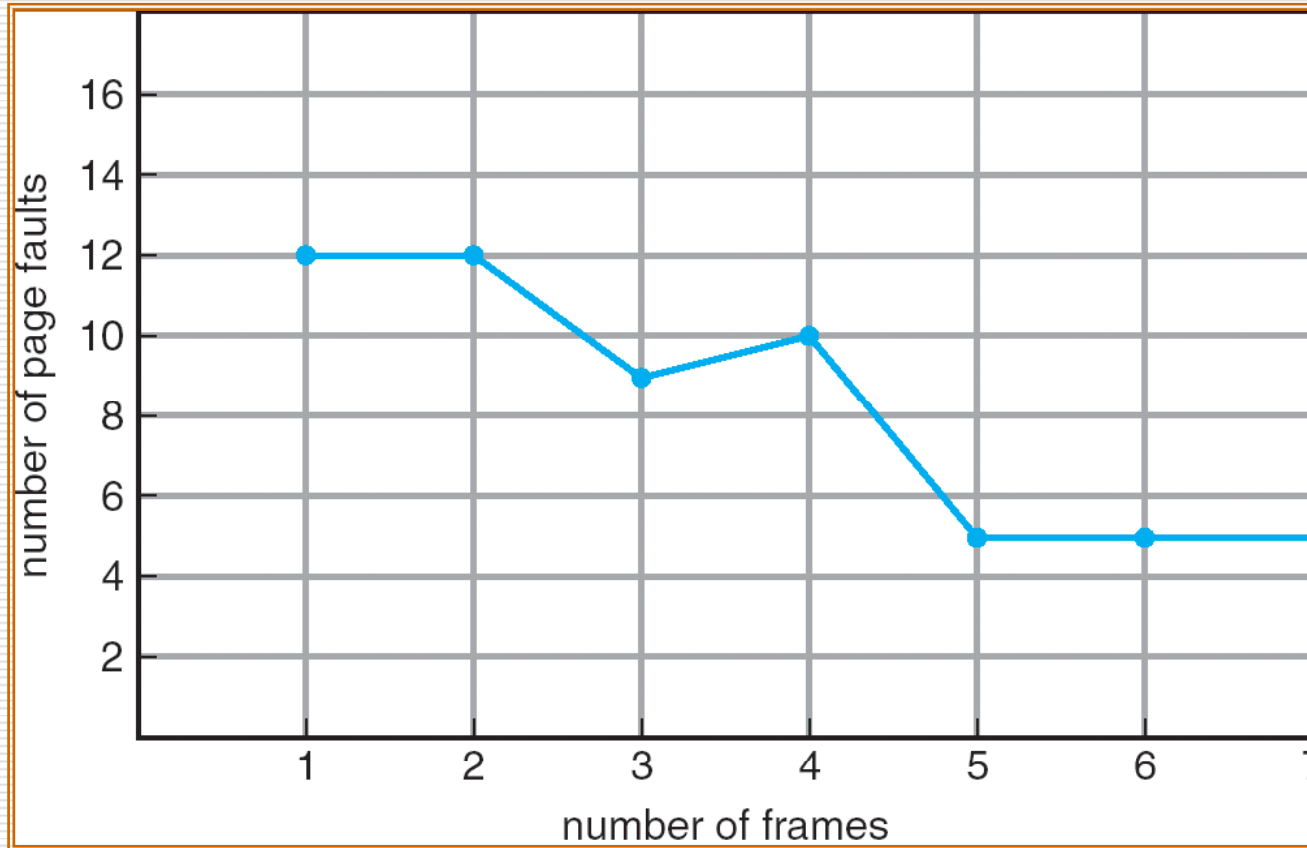
Belady's Anomaly

1 2 3 4 1 2 5 1 2 3 4 5



10 page faults

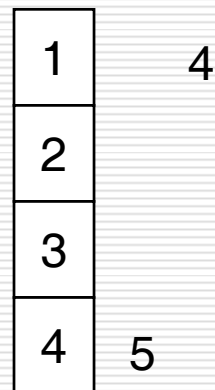
FIFO Illustrating Belady's Anomaly



Optimal Algorithm

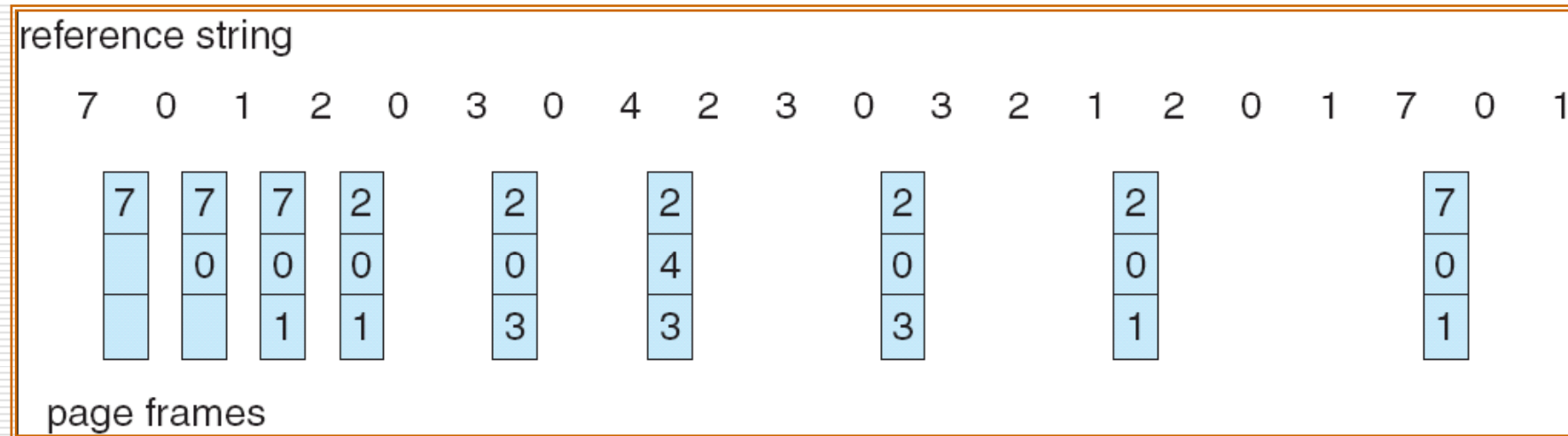
- ❑ Replace page that will not be used for longest period of time
- ❑ 4 frames example

1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5



6 page faults

Optimal Page Replacement

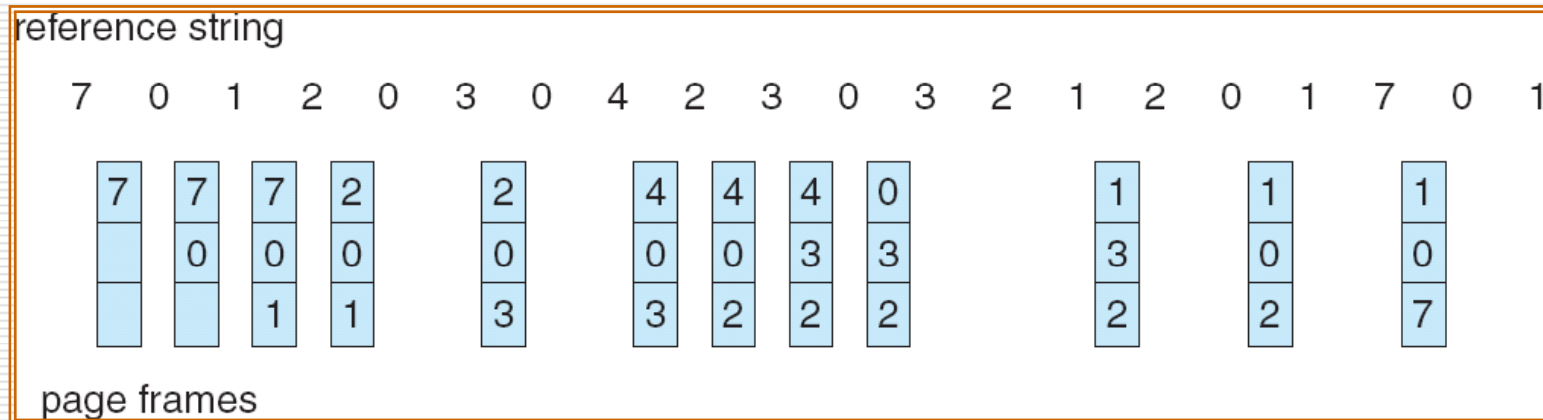


- How do you know this?
- Used for measuring how well your algorithm performs

Least Recently Used (LRU) Algorithm

□ Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

1	1	1	1	5
2	2	2	2	2
3	5	5	4	4
4	4	3	3	3



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 determine which are to change

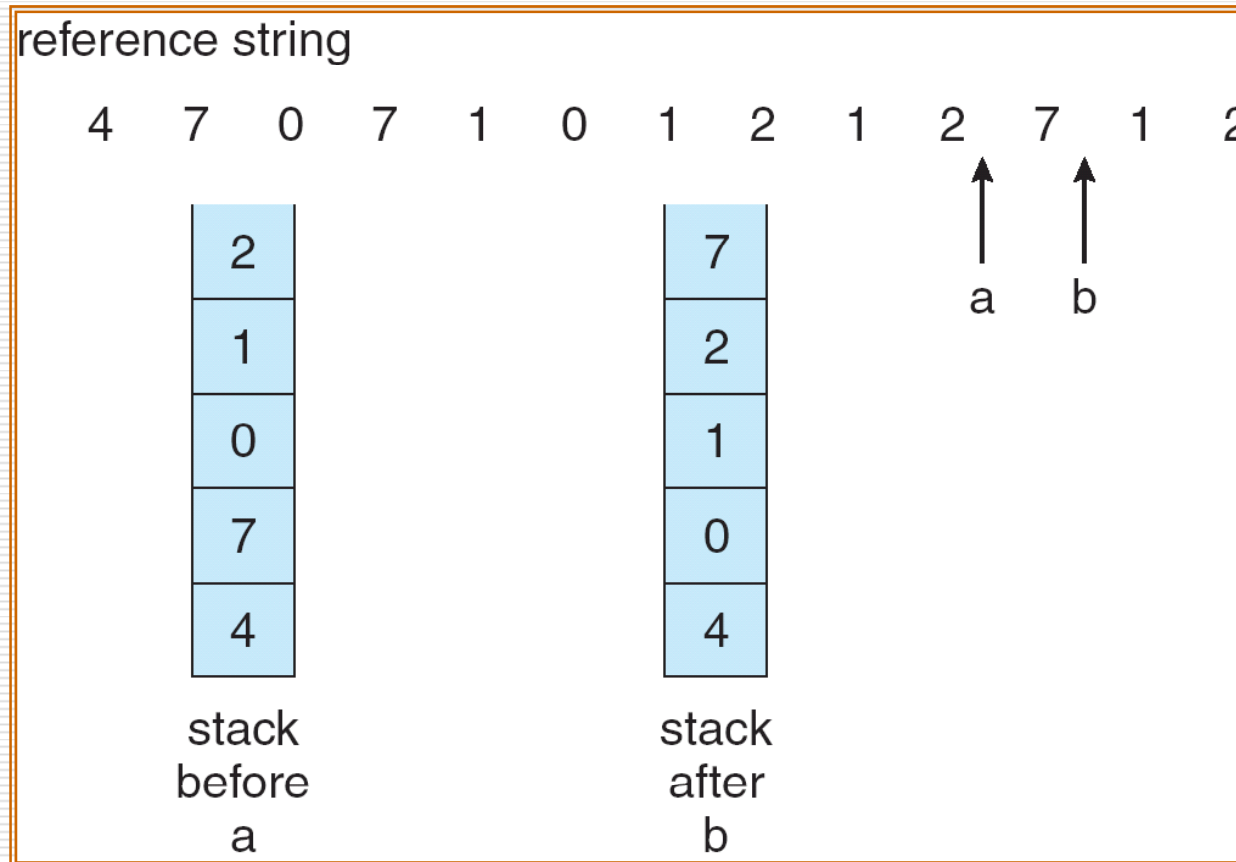
□ Problems

- Search time
- Overflow of the clock

LRU Algorithm (Cont.)

- ❑ 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
 - No search for replacement
 - ❑ The top of the stack is always the most recently used page
 - ❑ The bottom is the LRU page

Use Of A Stack to Record The Most Recent Page References



LRU Approximation Algorithms

□ Reference bit

- With each page associate a bit, initially = 0
- When page is referenced bit set to 1
- Replace the one which is 0 (if one exists)
 - We do not know the order, however

□ Additional-reference-bits algorithm

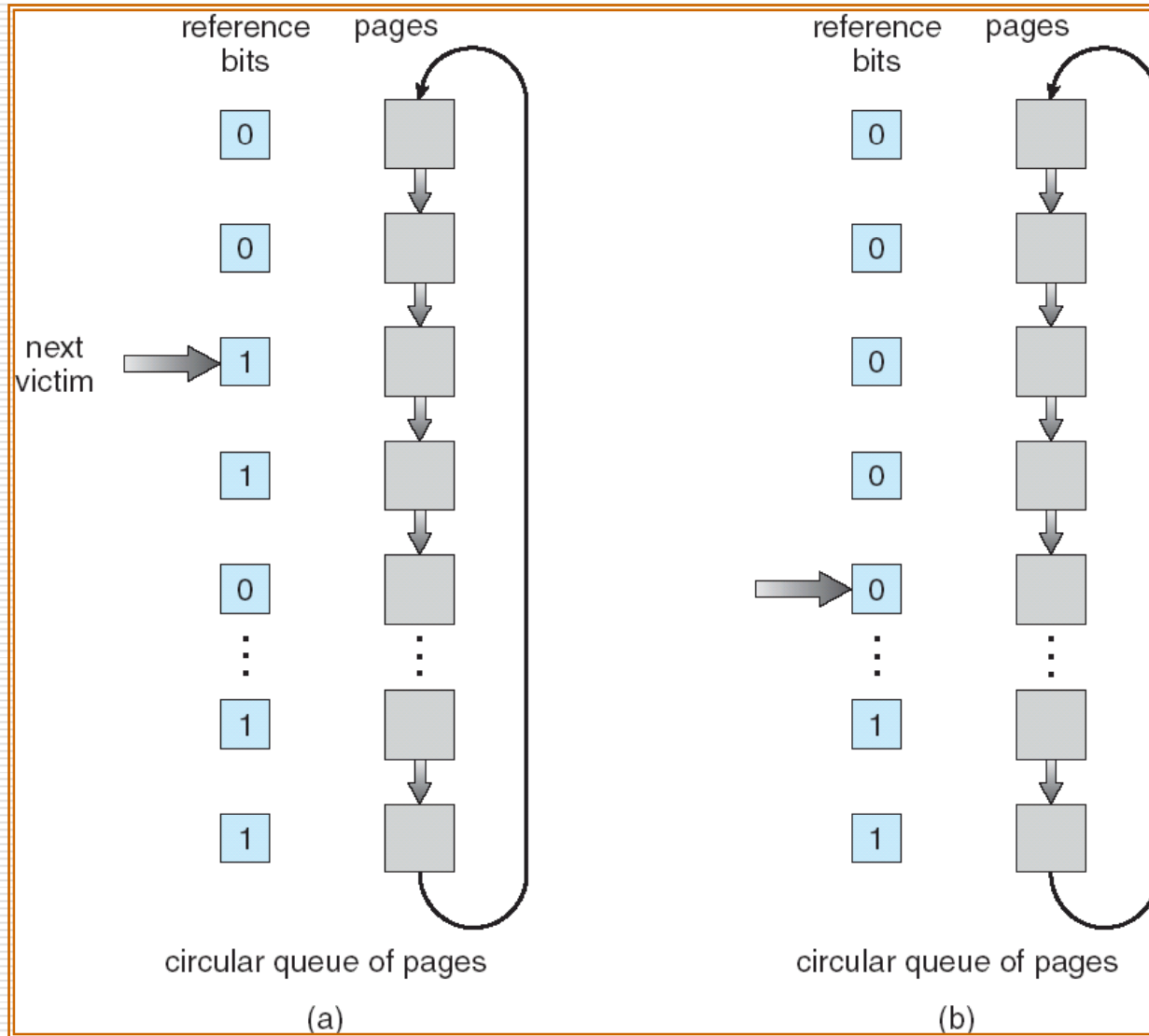
- Gain additional ordering information by recording the reference bits at regular intervals.
- 8bits
- Every 100 milliseconds

LRU Approximation Algorithms

Second chance

- Need reference bit
- Clock replacement
- If page to be replaced (in clock order) has reference bit = 1 then:
 - set reference bit 0
 - leave page in memory
 - replace next page (in clock order), subject to same rules

Second-Chance (clock) Page-Replacement Algorithm



Enhanced Second-Chance Algorithm

- Use an order pair: **Reference bit** and **Modify bit**
- 4 possible classes:
 - **(0,0)** neither recently used nor modified—best page to replace
 - **(0,1)** not recently used but modified—not quite as good. The page will be written out before replacement
 - **(1,0)** recently used but clean—it probably will be used again soon
 - **(1,1)** recently used and modified— **it probably will be used again soon**, and the page will be need to be written out to disk before it can be replaced
- Replace the first page encountered in the lowest nonempty class
- Compared clock algorithm, this algorithm give preference to those pages that have not been modified to reduce the number of I/Os required

Counting Algorithms

- Keep a counter of the number of references that have been made to each page
- **LFU Algorithm(Least Frequently Used):**
replaces page with smallest count
- **MFU Algorithm(Most Frequently Used):**
based on the argument that the page with the smallest count was probably just brought in and has yet to be used

Page-buffering algorithms

- ❑ Keep a pool of free frames
- ❑ The desired page is read into a free frame from the pool before the victim is written out.
- ❑ When the victim is later written out, its frame is added to the free-frame pool.

- ❑ This method can be used combined with other algorithms.

Allocation of Frames

- Each process needs *minimum* number of pages
- Example: IBM 370 – 6 pages to handle SS MOVE instruction:
 - instruction is 6 bytes, might span 2 pages
 - 2 pages to handle *from*
 - 2 pages to handle *to*
- Two major allocation schemes
 - fixed allocation
 - priority allocation

Fixed Allocation

- Equal allocation – For example, if there are 100 frames and 5 processes, give each process 20 frames.
- Proportional allocation – Allocate according to the size of process

s_i = size of process p_i

$$S = \sum s_i$$

m = total number of frames

$$a_i = \text{allocation for } p_i = \frac{s_i}{S} \times m$$

$$m = 64$$

$$s_1 = 10$$

$$s_2 = 127$$

$$a_1 = \frac{10}{137} \times 64 \approx 5$$

$$a_2 = \frac{127}{137} \times 64 \approx 59$$

Fixed Allocation

□ 特点

- 每个进程所分配的数量会随着多道程序的级别而有所变化。多道程序程度增加，那么每个进程会失去一些帧以提供给新来进程使用。反之，原来分配给离开进程的帧可以分配给剩余进程
- 高优先级进程与低优先级进程在这种分配方式下没有任何区别

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

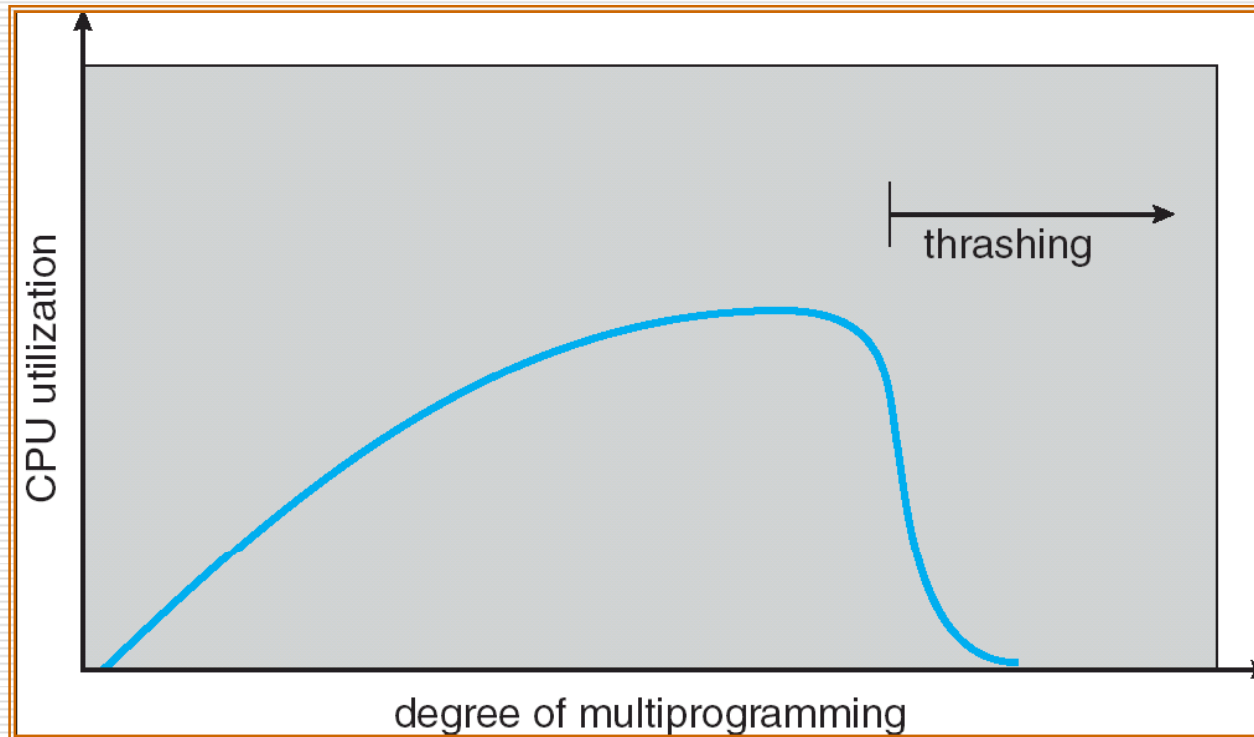
Global vs. Local Allocation

- ❑ **Global replacement** – process selects a replacement frame from the set of all frames; one process can take a frame from another
- ❑ **Local replacement** – each process selects from only its own set of allocated frames

Thrashing

- 颠簸：进程的频繁的页调度行为
 - 页错误显著增加
 - 吞吐量徒降
 - 有效访问时间增加
- If a process does not have “enough” pages, the page-fault rate is very high. This leads to:
 - low CPU utilization
 - operating system thinks 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?

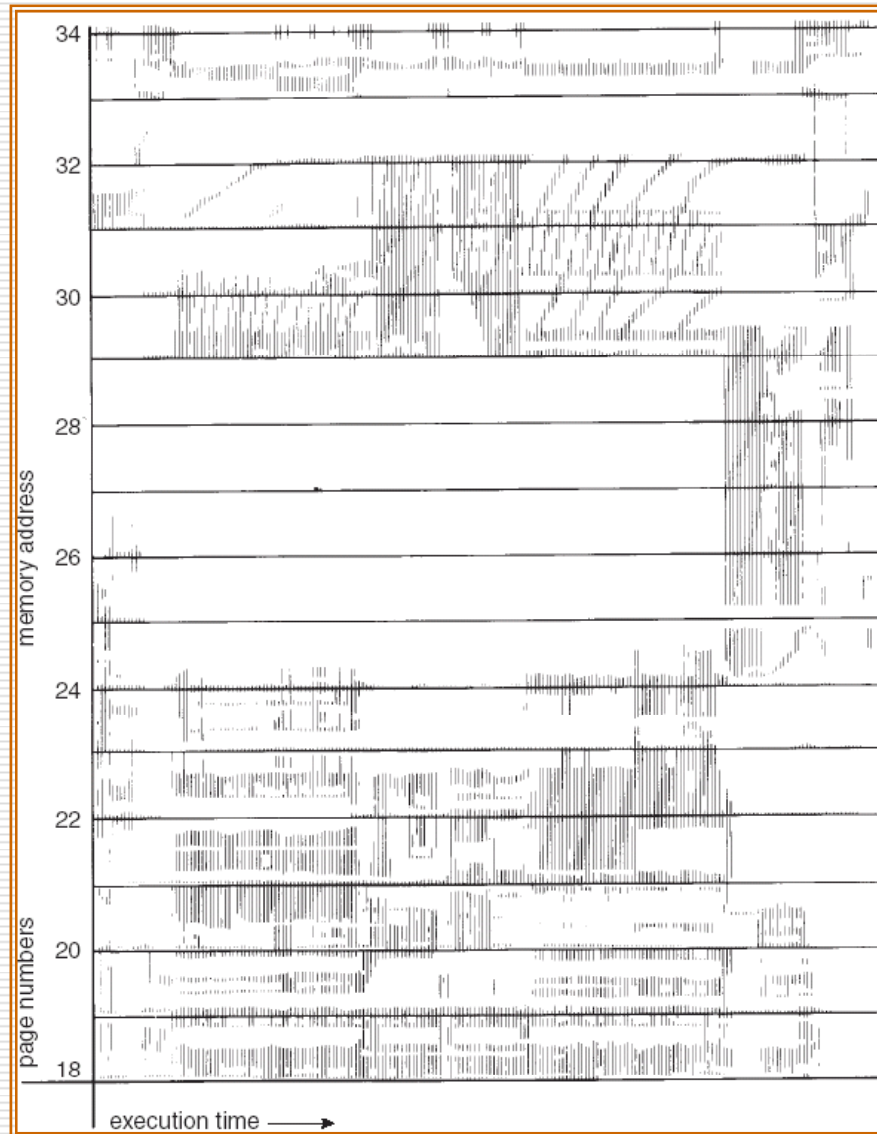
Locality model

- Process migrates from one locality to another
- Localities may overlap

□ Why does thrashing occur?

Σ size of locality > total memory size

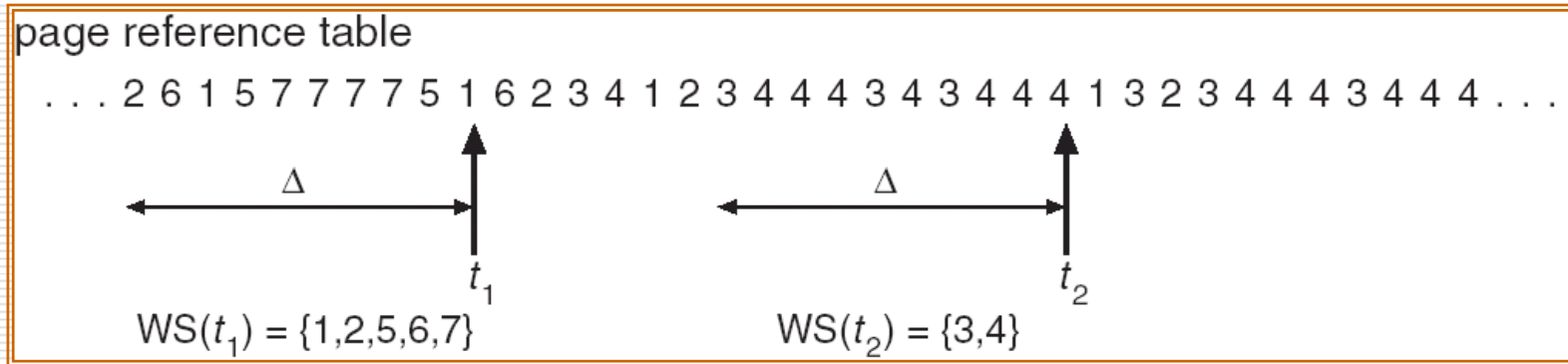
Locality In A Memory-Reference Pattern



Working-Set Model

- $\Delta \equiv$ working-set window \equiv a fixed number of page references
Example: 10,000 instruction
- WSS_i (working set of Process P_i) = total number of pages referenced in the most recent Δ (varies in time)
 - if Δ too small will not encompass entire locality
 - if Δ too large will encompass several localities
 - if $\Delta = \infty \Rightarrow$ will encompass entire program
- $D = \sum WSS_i \equiv$ total demand frames
- if $D > m \Rightarrow$ Thrashing
- Policy if $D > m$, then suspend one of the processes

Working-set model

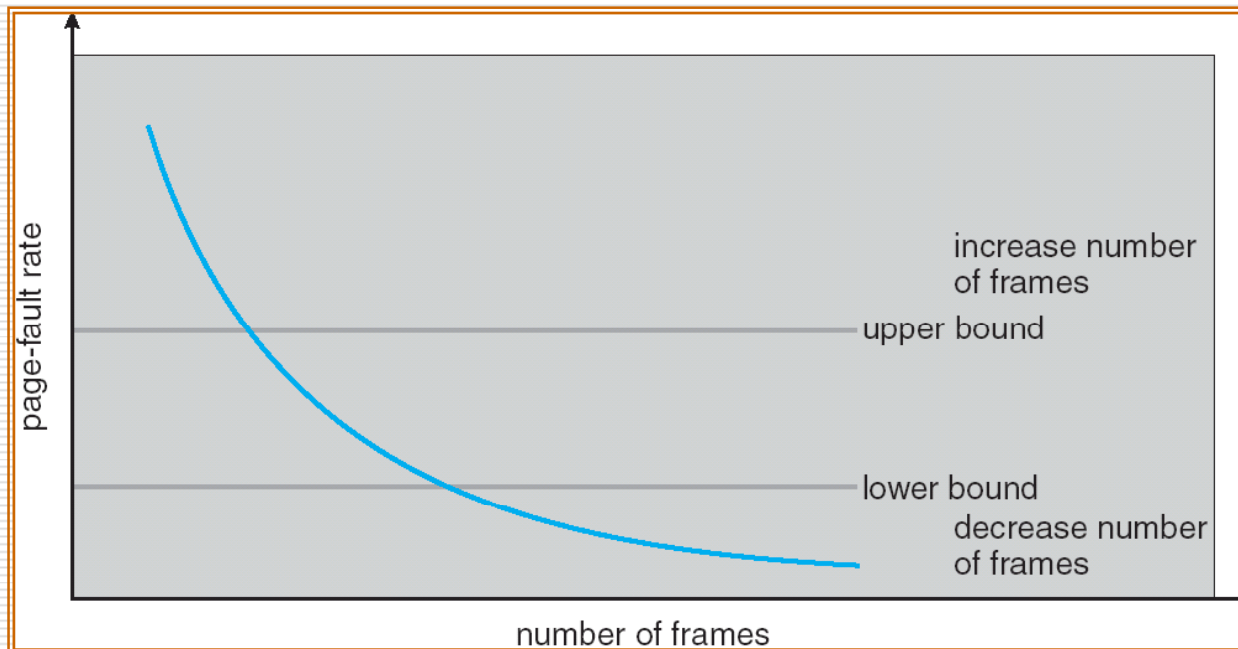


Keeping Track of the Working Set

- Approximate with interval timer + a reference bit
- Example: $\Delta = 10,000$
 - Timer interrupts after every 5000 time units
 - Keep in memory 2 bits for each page
 - Whenever a timer interrupts copy and set the values of all reference bits to 0
 - If one of the bits in memory = 1 \Rightarrow page in working set
- Why is this not completely accurate?
- Improvement = 10 bits and interrupt every 1000 time units

Page-Fault Frequency Scheme

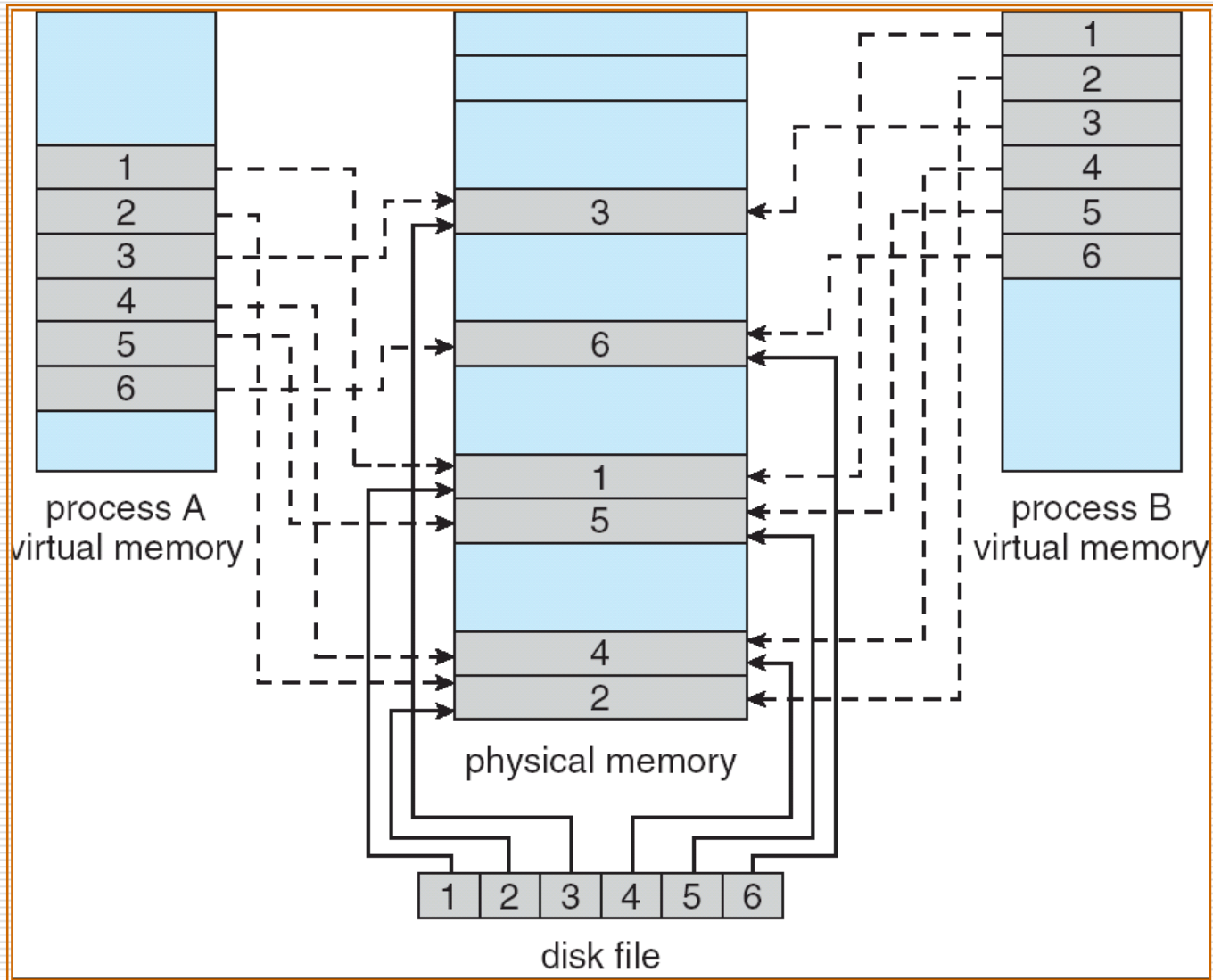
- Establish “acceptable” page-fault rate
 - If actual rate too low, process loses frame
 - If actual rate too high, process gains frame



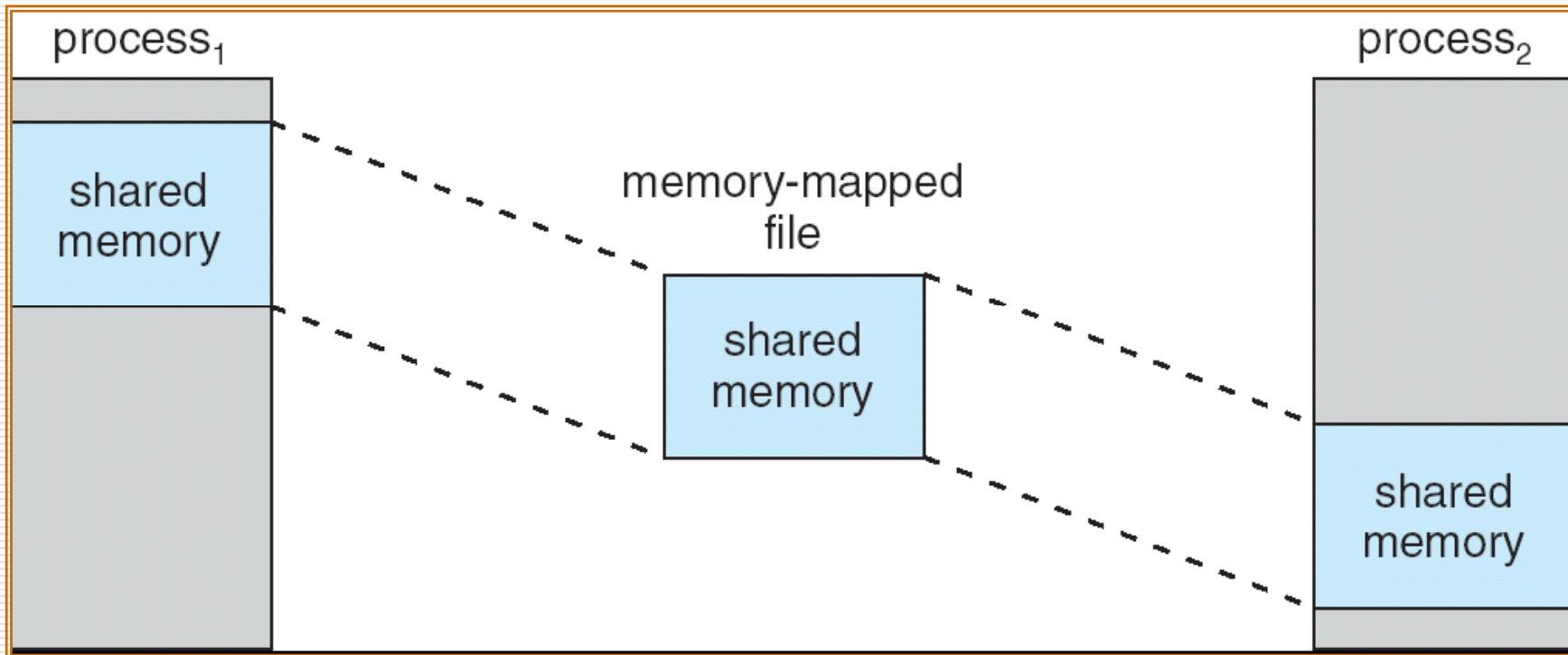
Memory-Mapped Files

- ❑ Memory-mapped file I/O allows file I/O to be treated as routine memory access by **mapping** a disk block to a page in memory
- ❑ A file is initially read using demand paging. A page-sized portion of the file is read from the file system into a physical page. Subsequent reads/writes to/from the file are treated as ordinary memory accesses.
- ❑ Simplifies file access by treating file I/O through memory rather than **read() write()** system calls
- ❑ Also allows several processes to map the same file allowing the pages in memory to be shared

Memory Mapped Files



Memory-Mapped Shared Memory in Windows



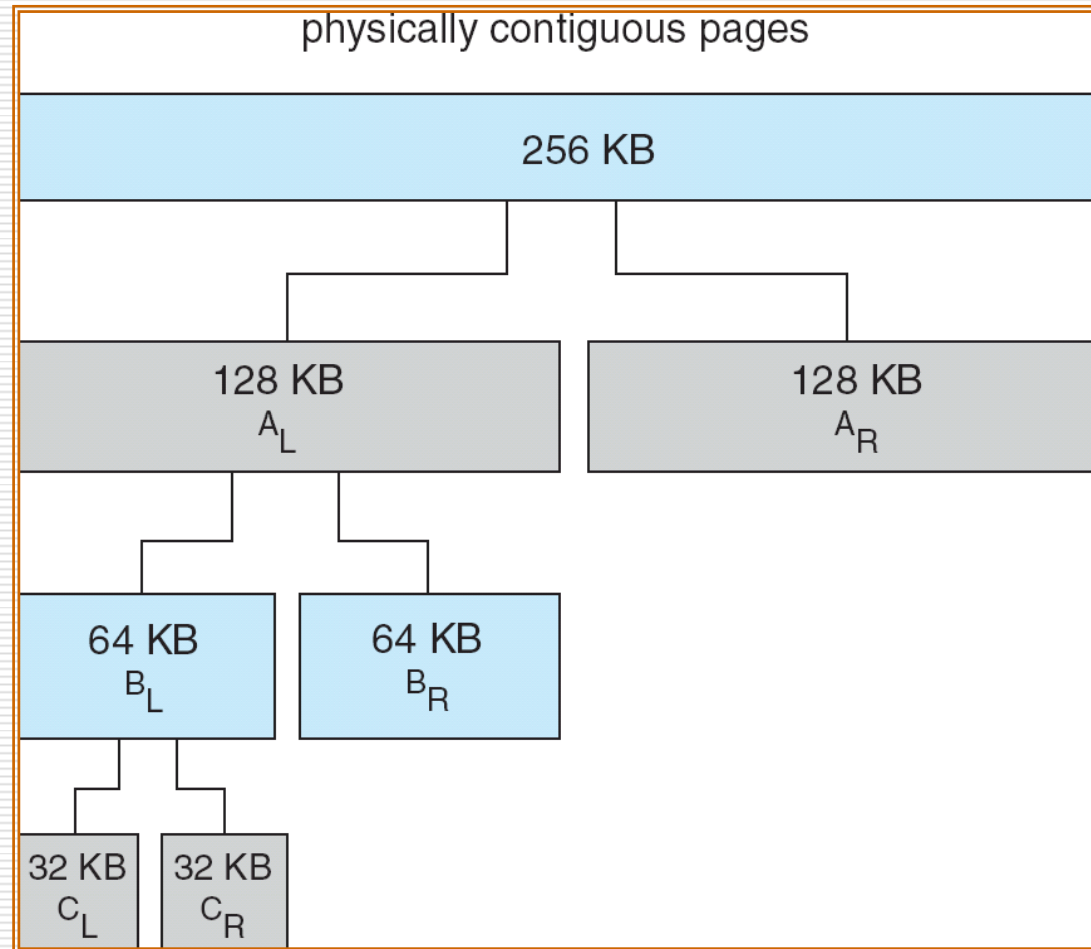
Allocating Kernel Memory

- ❑ 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

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

Buddy System Allocator



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

Slab Class: 1

Chunks:

88 bytes

88 bytes

88 bytes

... lots more!

Slab Class: 2

Chunks:

112 bytes

112 bytes

112 bytes

... lots more!

Slab Class: 3

Chunks:

144 bytes

144 bytes

144 bytes

... lots more!

Slab Class: n

Chunks:

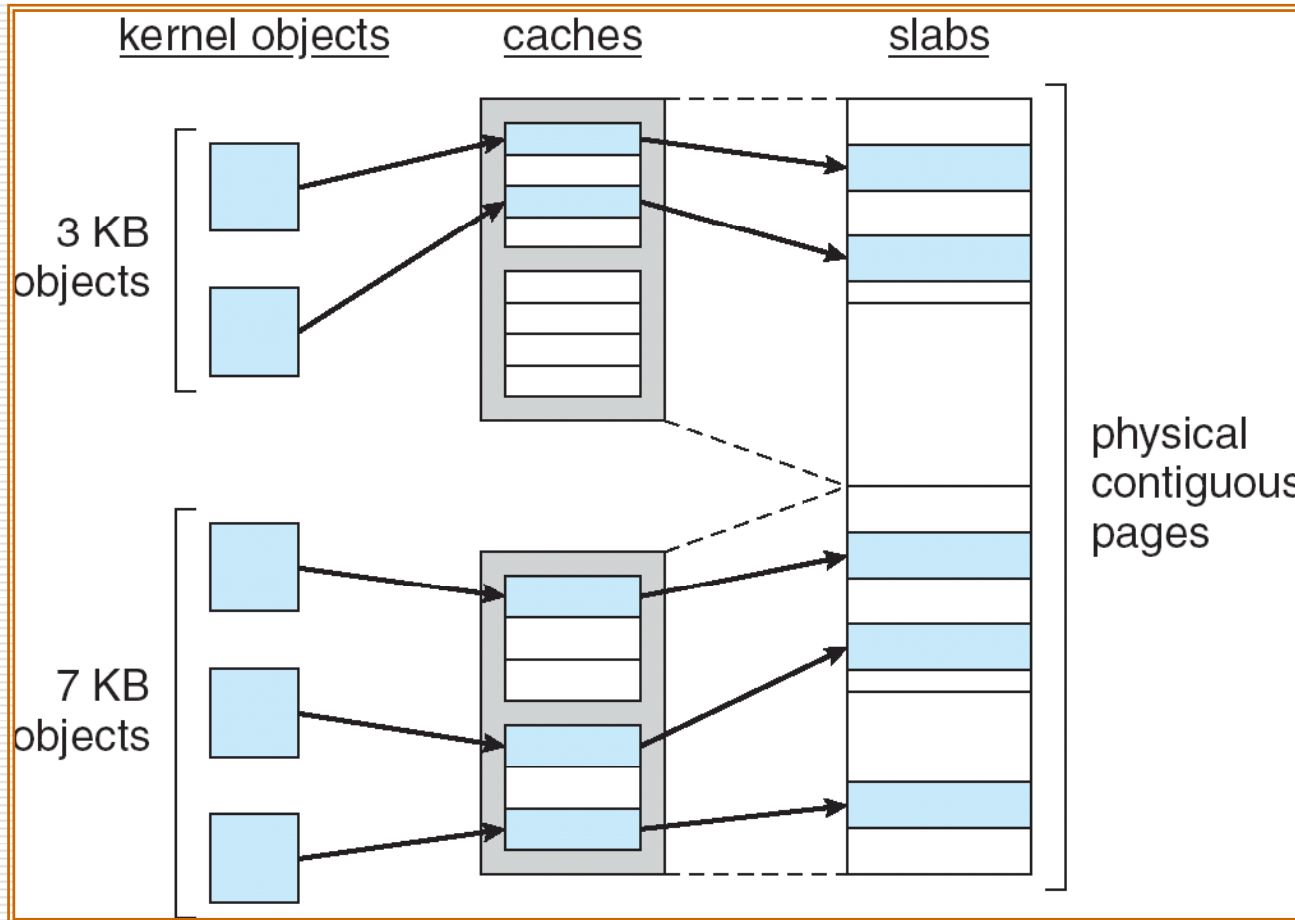
n bytes

n bytes

n bytes

... more!

Slab Allocation



Other Issues -- Prepaging

□ Prepaging

- To reduce the large number of page faults that occurs at process startup
- Prepage all or some of the pages a process will need, before they are referenced
- But if prepaged pages are unused, I/O and memory was wasted
- Assume s pages are prepaged and α of the pages is used
 - Is cost of $s * \alpha$ save pages faults $>$ or $<$ than the cost of prepaging $s * (1 - \alpha)$ unnecessary pages?
 - α near zero \Rightarrow prepaging loses

Other Issues – Page Size

- Page size selection must take into consideration:
 - fragmentation
 - table size
 - I/O overhead
 - locality

Other Issues – TLB Reach

- ❑ TLB Reach - The amount of memory accessible from the TLB
- ❑ TLB Reach = (TLB Size) X (Page Size)
- ❑ Ideally, the working set of each process is stored in the TLB
 - Otherwise there is a high degree of page faults
- ❑ Increase the Page Size
 - This may lead to an increase in fragmentation as not all applications require a large page size
- ❑ Provide Multiple Page Sizes
 - This allows applications that require larger page sizes the opportunity to use them without an increase in fragmentation

Other Issues – Program Structure

- Program structure

- Int[128,128] data;
- Each row is stored in one page
- Program 1

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

128 x 128 = 16,384 page faults

- Program 2

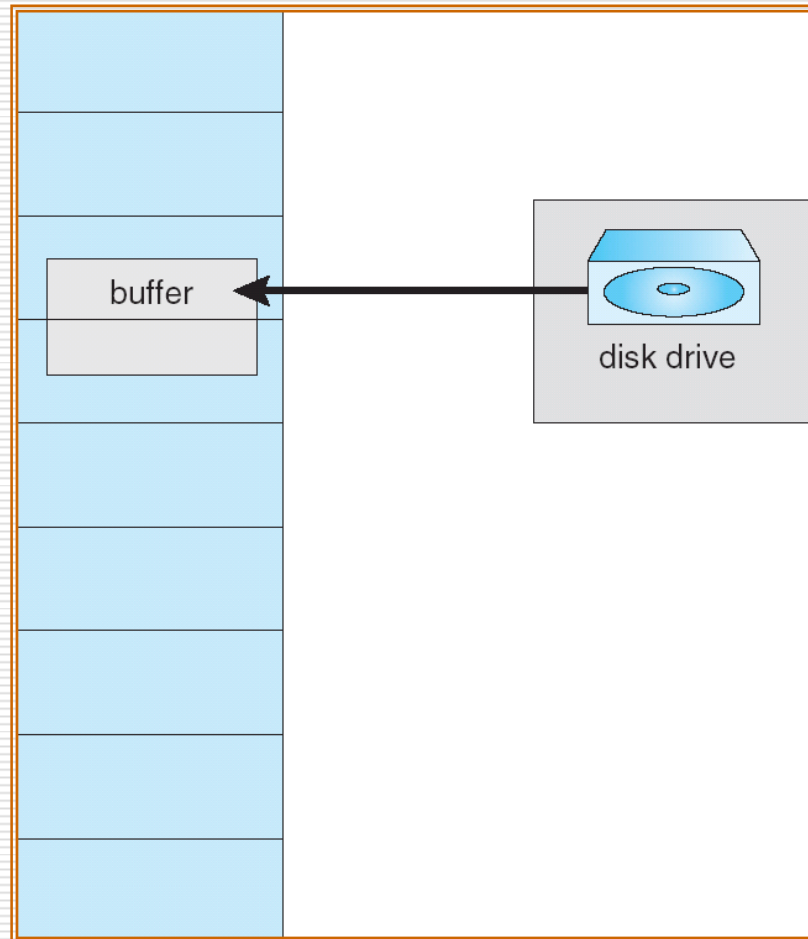
```
for (i = 0; i < 128; i++)  
    for (j = 0; j < 128; j++)  
        data[i,j] = 0;
```

128 page faults

Other Issues – I/O interlock

- ❑ **I/O Interlock** – Pages must sometimes be locked into memory
- ❑ Consider I/O - Pages that are used for copying a file from a device must be locked from being selected for eviction by a page replacement algorithm

Reason Why Frames Used For I/O Must Be In Memory



Operating System Examples

Windows XP

Solaris

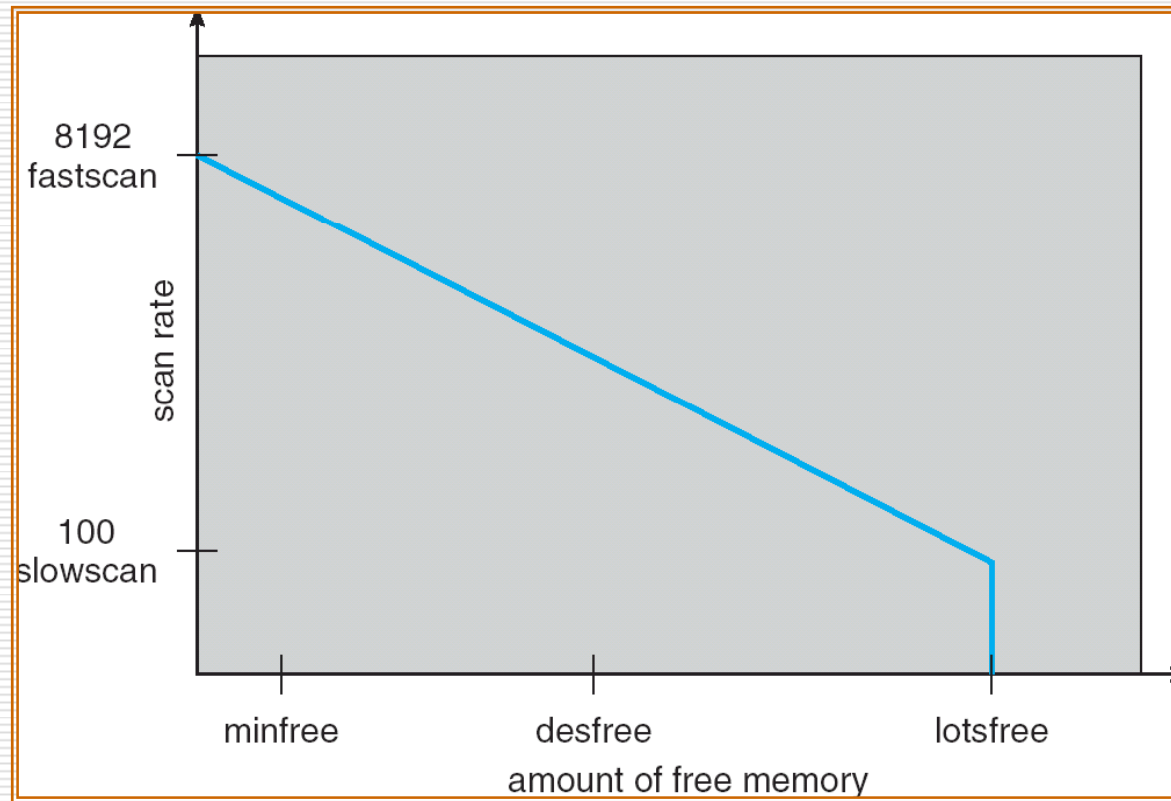
Windows XP

- ❑ Uses demand paging with **clustering**. Clustering brings in pages surrounding the faulting page.
- ❑ Processes are assigned **working set minimum** and **working set maximum**
- ❑ Working set minimum is the minimum number of pages the process is guaranteed to have in memory
- ❑ A process may be assigned as many pages up to its working set maximum
- ❑ When the amount of free memory in the system falls below a threshold, **automatic working set trimming** is performed to restore the amount of free memory
- ❑ Working set trimming removes pages from processes that have pages in excess of their working set minimum

Solaris

- ❑ Maintains a list of free pages to assign faulting processes
- ❑ *Lotsfree* – threshold parameter (amount of free memory) to begin paging
- ❑ *Desfree* – threshold parameter to increasing paging
- ❑ *Minfree* – threshold parameter to being swapping
- ❑ Paging is performed by *pageout* process
- ❑ Pageout scans pages using modified clock algorithm
- ❑ *Scanrate* is the rate at which pages are scanned. This ranges from *slowscan* to *fastscan*
- ❑ Pageout is called more frequently depending upon the amount of free memory available

Solaris 2 Page Scanner



Assignments

□ 9.4 9.10 9.17 9.18

End of Chapter 9

Any Question?