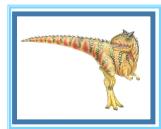


Chapter 9: Virtual Memory



Operating System Concepts – 8th Edition,

Silberschatz, Galvin and Gagne ©2009



Background: True or False

- Instructions being executed (and data being used) must be in physical memory (RAM)
- All instructions and all data must always be in RAM



Silberschatz, Galvin and Gagne ©2009

Operating System

9.2



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



Operating System Concepts – 8th Edition

9.3

Silberschatz, Galvin and Gagne ©2009



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



Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition

9.4



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 can be implemented via:
 - Demand paging
 - Demand segmentation



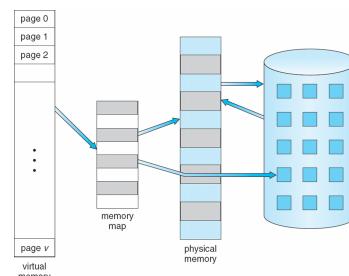
Operating System Concepts – 8th Edition

9.5

Silberschatz, Galvin and Gagne ©2009



Virtual Memory That is Larger Than Physical Memory

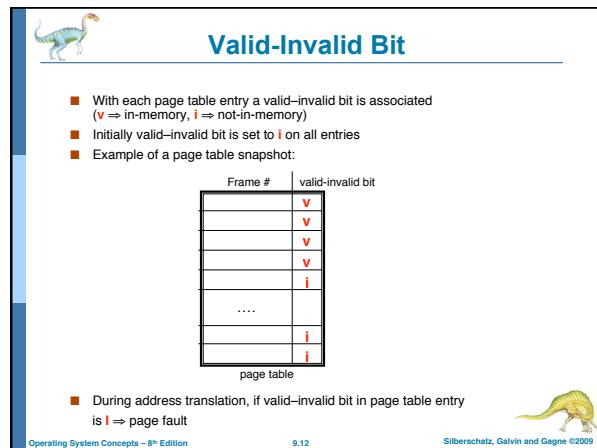
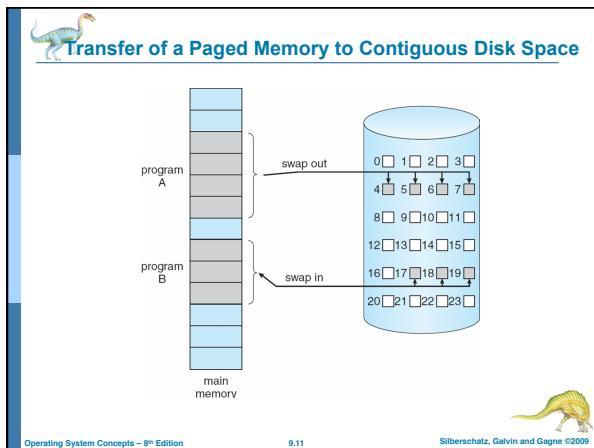
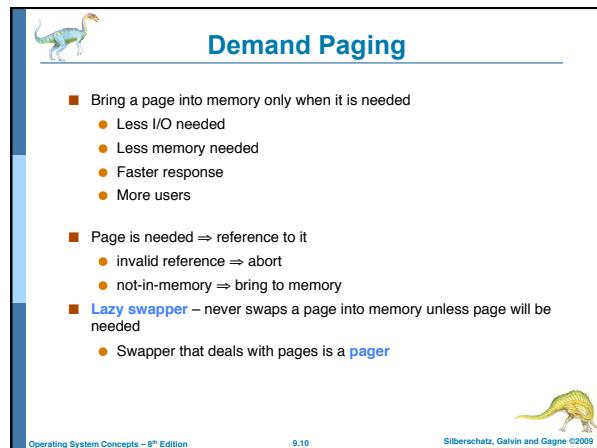
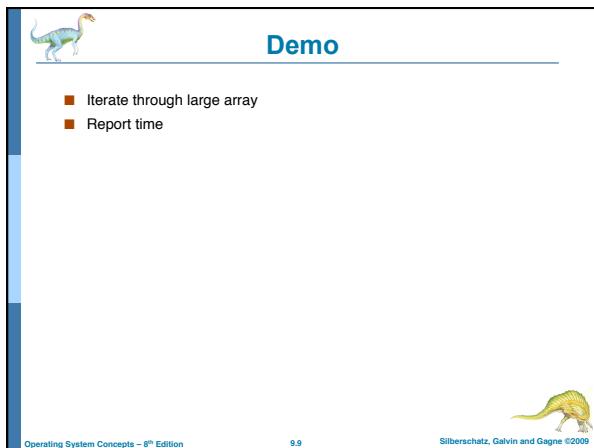
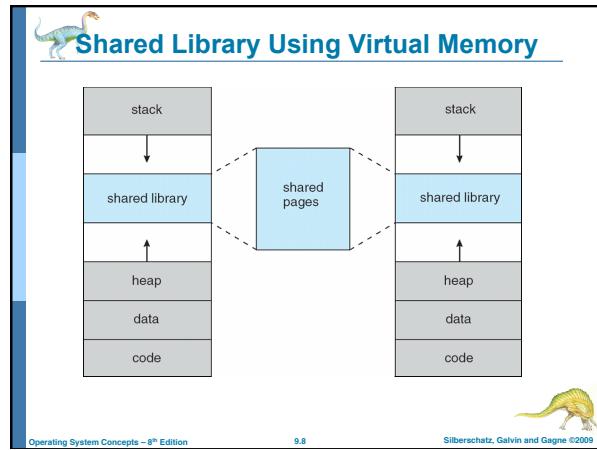
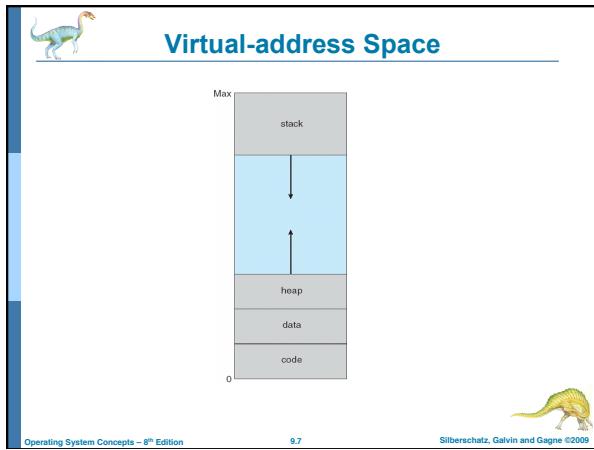


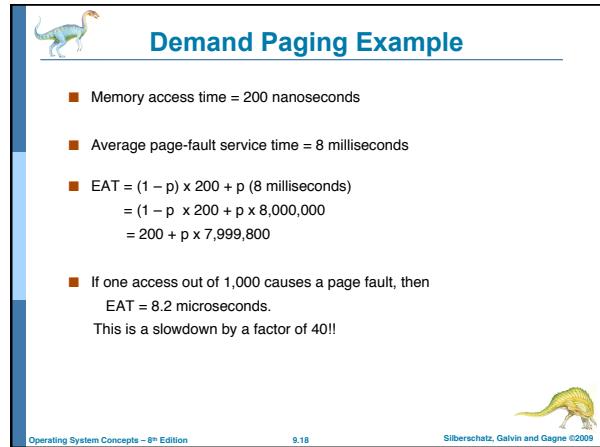
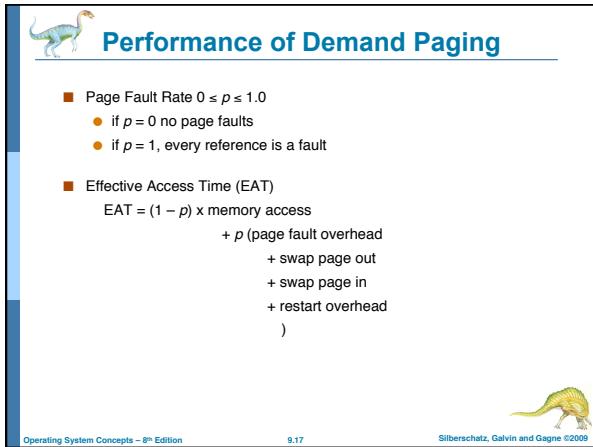
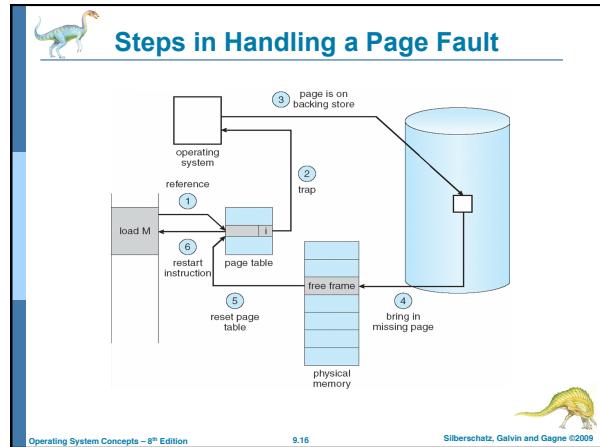
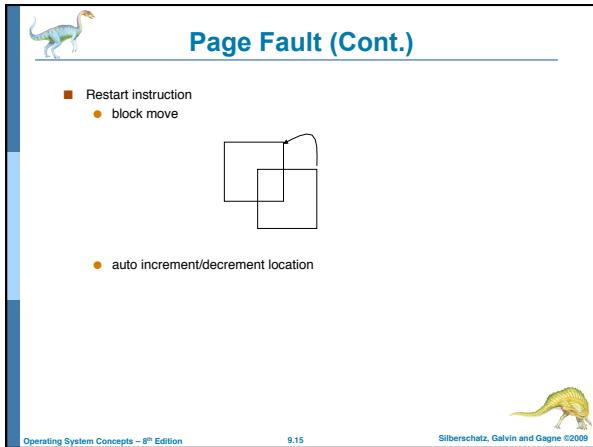
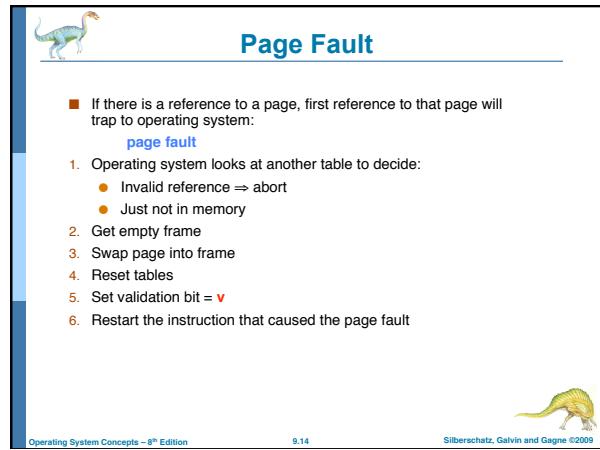
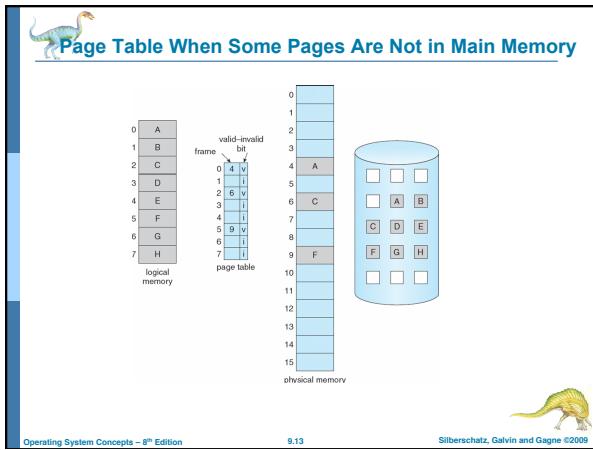
Operating System Concepts – 8th Edition

9.6

Silberschatz, Galvin and Gagne ©2009









Process Creation

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

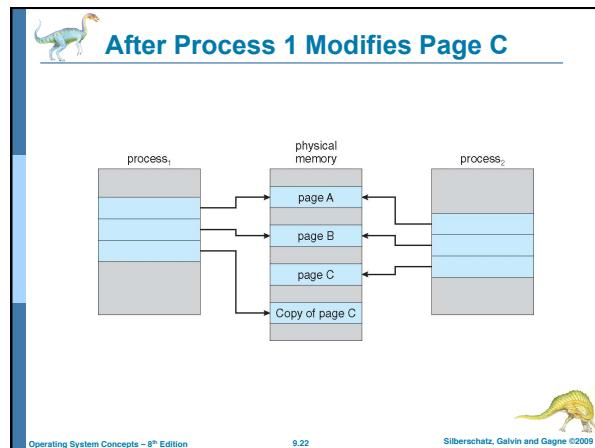
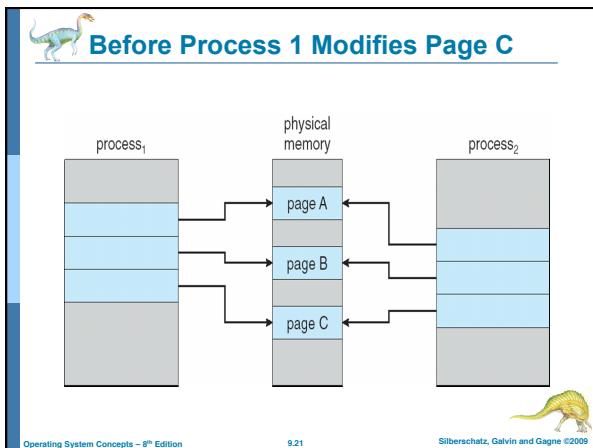
Operating System Concepts – 8th Edition 9.19 Silberschatz, Galvin and Gagne ©2009



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

Operating System Concepts – 8th Edition 9.20 Silberschatz, Galvin and Gagne ©2009




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

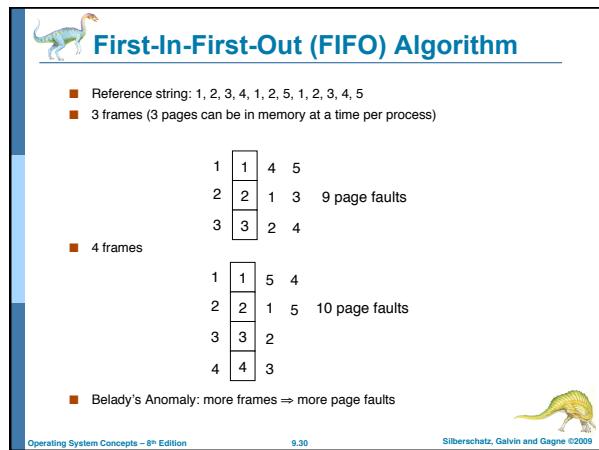
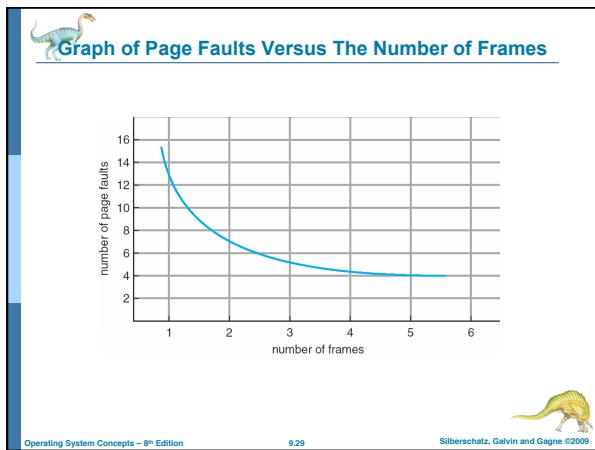
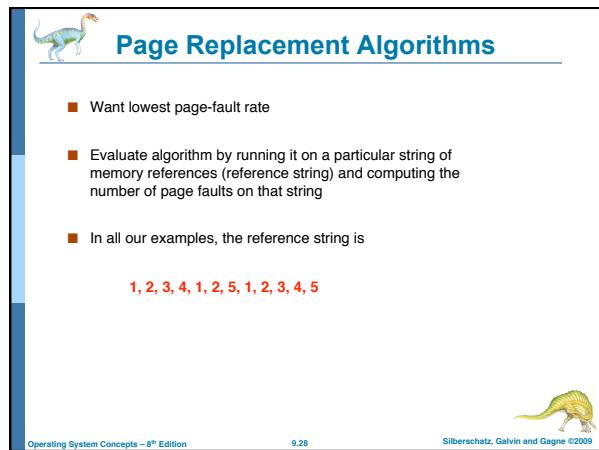
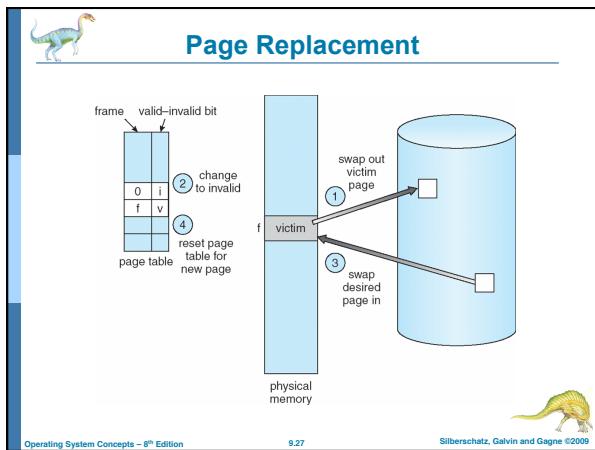
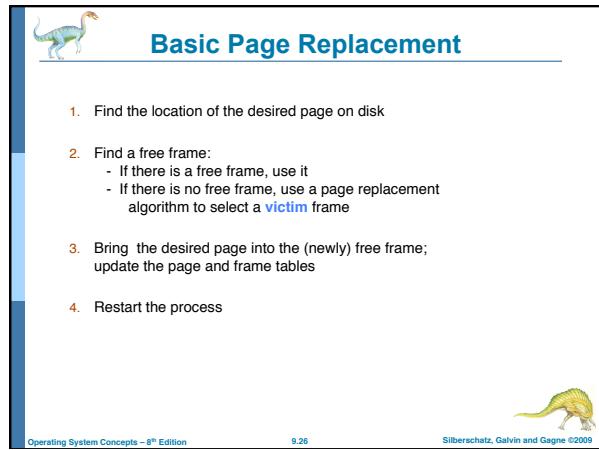
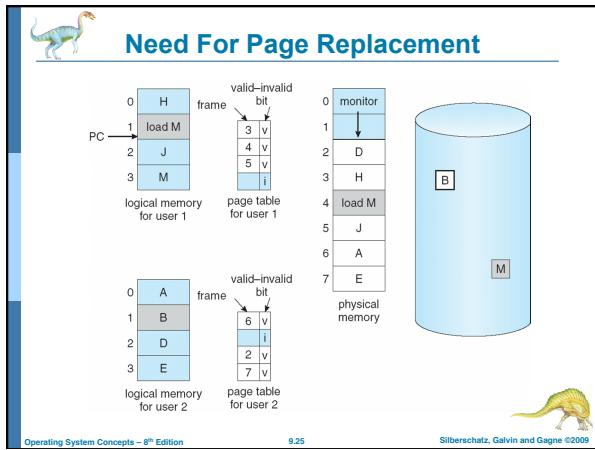
Operating System Concepts – 8th Edition 9.23 Silberschatz, Galvin and Gagne ©2009

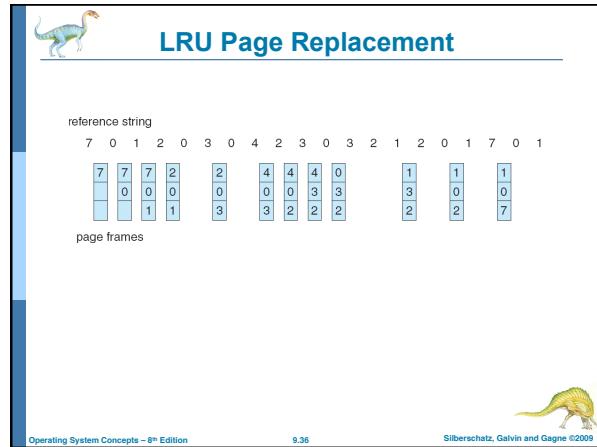
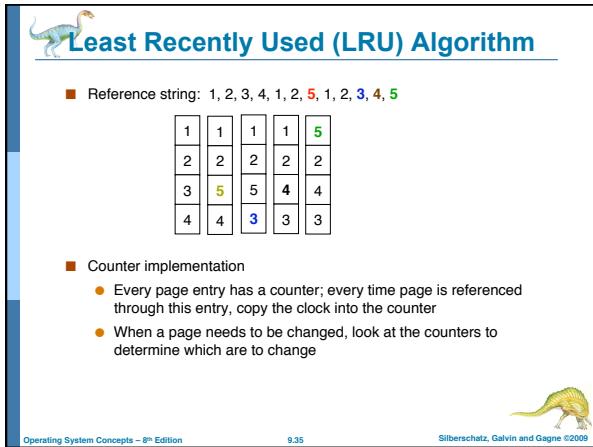
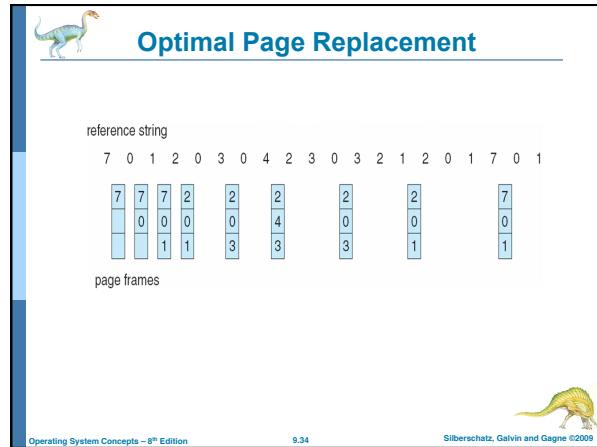
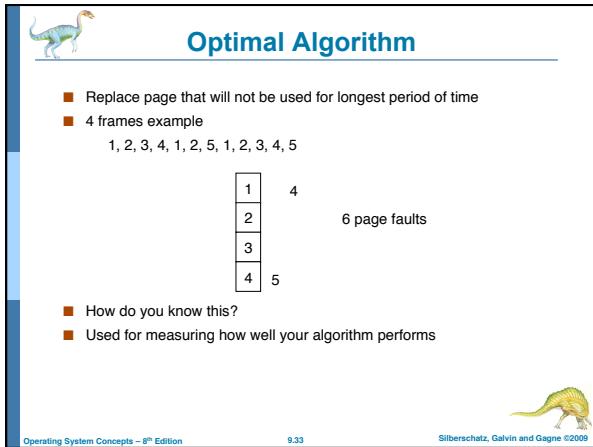
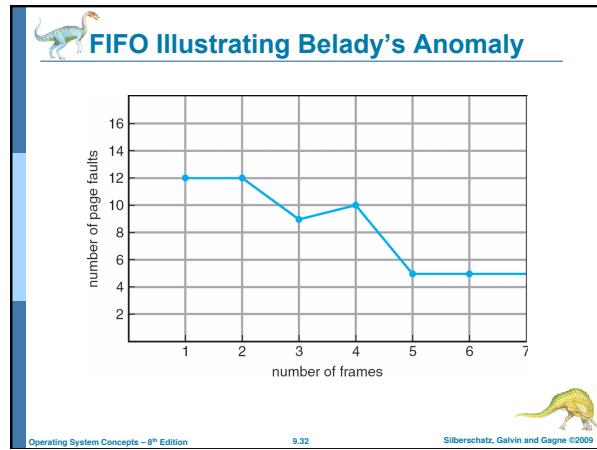
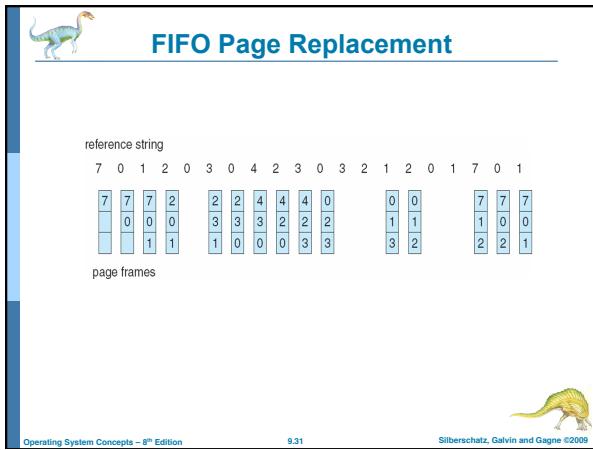


Page Replacement

- 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

Operating System Concepts – 8th Edition 9.24 Silberschatz, Galvin and Gagne ©2009



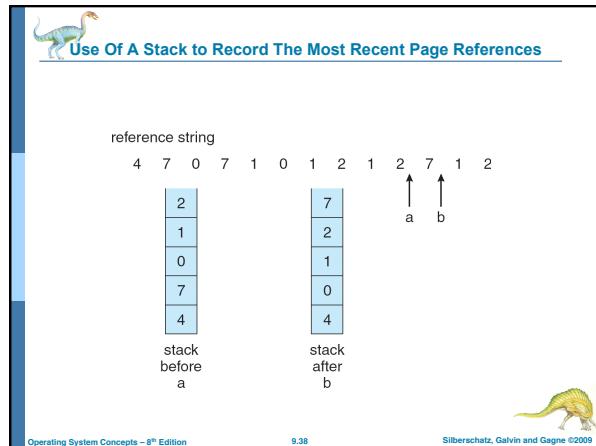




LRU Algorithm (Cont.)

- 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

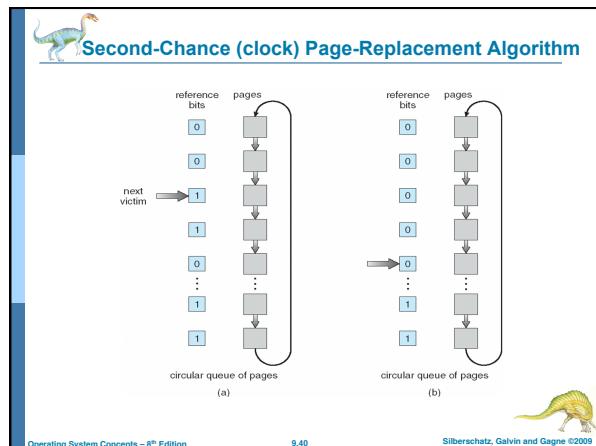
Operating System Concepts – 8th Edition 9.37 Silberschatz, Galvin and Gagne ©2009




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
- 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

Operating System Concepts – 8th Edition 8.39 Silberschatz, Galvin and Gagne ©2009




Counting Algorithms

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

Operating System Concepts – 8th Edition 9.41 Silberschatz, Galvin and Gagne ©2009



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

Operating System Concepts – 8th Edition 9.42 Silberschatz, Galvin and Gagne ©2009



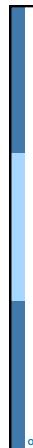
Fixed Allocation

- 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 = allocation for $p_i = \frac{S_i}{S} \times m$

$m = 64$
 $S_1 = 10$
 $S_2 = 127$
 $a_1 = \frac{10}{137} \times 64 \approx 5$
 $a_2 = \frac{127}{137} \times 64 \approx 59$


Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition 9.43 Silberschatz, Galvin and Gagne ©2009



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


Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition 9.44 Silberschatz, Galvin and Gagne ©2009

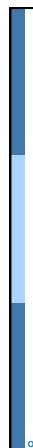


Global vs. Local Allocation

- 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


Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition 9.45 Silberschatz, Galvin and Gagne ©2009



Demonstration

NAME
vmstat - report virtual memory statistics

SYNOPSIS
vmstat [-cipsS] [disks] [interval [count]]

DESCRIPTION
vmstat reports virtual memory statistics regarding process, virtual memory, disk, trap, and CPU activity.


Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition 9.46 Silberschatz, Galvin and Gagne ©2009



vmstat

The fields of vmstat's display are

procs Report the number of processes in each of the three following states:

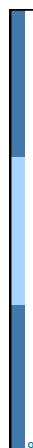
- r in run queue
- b blocked for resources I/O, paging, and so forth
- w runnable but swapped

memory Report on usage of virtual and real memory.

- swap amount of swap space currently available (Kbytes)
- free size of the free list (Kbytes)


Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition 9.47 Silberschatz, Galvin and Gagne ©2009



vmstat

page Report information about page faults and paging activity. The information on each of the following activities is given in units per second.

re page reclaims - but see the **-S** option for how this field is modified.

mf minor faults - but see the **-S** option for how this field is modified.

pi kilobytes paged in

po kilobytes paged out

fr kilobytes freed

de anticipated short-term memory shortfall (Kbytes)

sr pages scanned by clock algorithm


Silberschatz, Galvin and Gagne ©2009

Operating System Concepts – 8th Edition 9.48 Silberschatz, Galvin and Gagne ©2009

  **vmstat**

Maintenance Commands vmstat(1M)

disk Report the number of disk operations per second.
faults Report the trap/interrupt rates (per second).

in (non clock) device interrupts
sy system calls
cs CPU context switches

cpu Give a breakdown of percentage usage of CPU time. On MP systems, this is an average across all processors.

us user time
sy system time
id idle time

Operating System Concepts – 8th Edition 9.49 Silberschatz, Galvin and Gagne ©2009

  **Demonstration**

- Thrash1: allocate a big piece of memory
- Thrash2: allocate a big piece of memory and access it sequentially
- Thrashmem: allocate a big piece of memory and access it randomly
- Rss1: allocate a big piece of memory and sleep
- Rss2: allocate a big piece of memory and access it sequentially, then sleep
- Watch also process start up

Operating System Concepts – 8th Edition 9.50 Silberschatz, Galvin and Gagne ©2009

  **Comments**

- Allocating memory:
- Traversing memory:

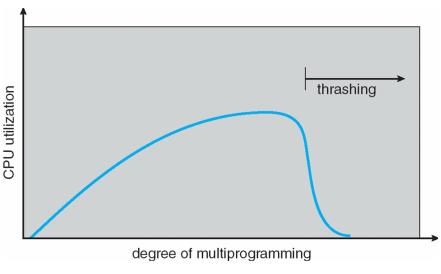
Operating System Concepts – 8th Edition 9.51 Silberschatz, Galvin and Gagne ©2009

  **Thrashing**

- 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** = a process is busy swapping pages in and out

Operating System Concepts – 8th Edition 9.52 Silberschatz, Galvin and Gagne ©2009

  **Thrashing (Cont.)**



CPU utilization

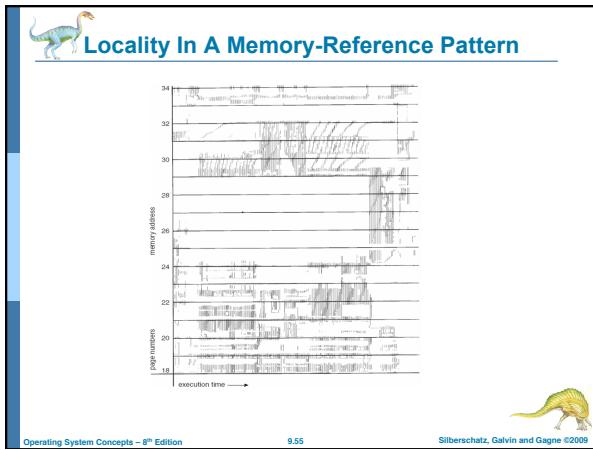
degree of multiprogramming

Operating System Concepts – 8th Edition 9.53 Silberschatz, Galvin and Gagne ©2009

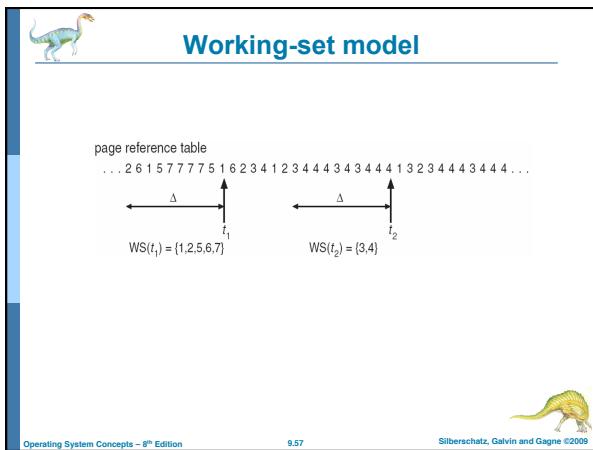
  **Demand Paging and Thrashing**

- 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

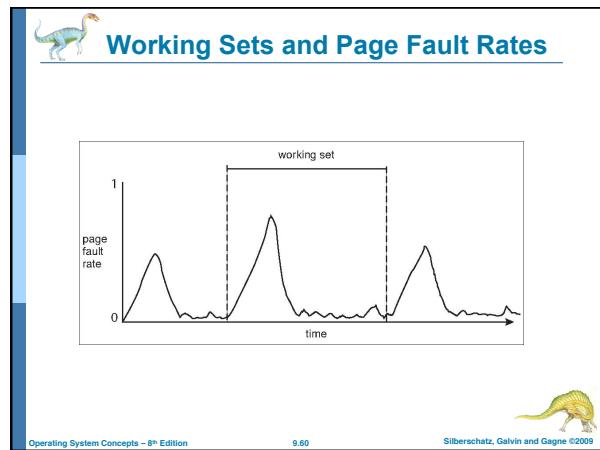
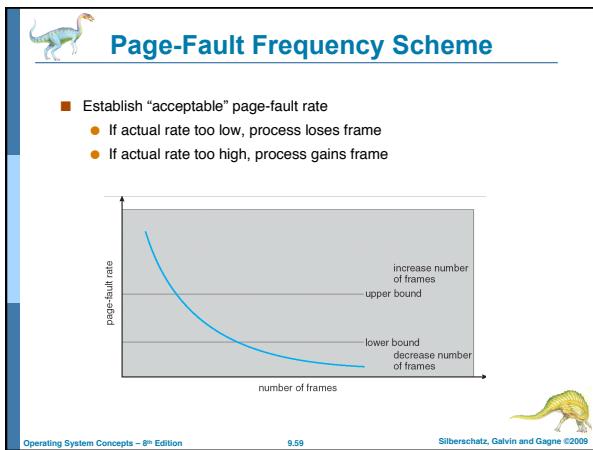
Operating System Concepts – 8th Edition 9.54 Silberschatz, Galvin and Gagne ©2009

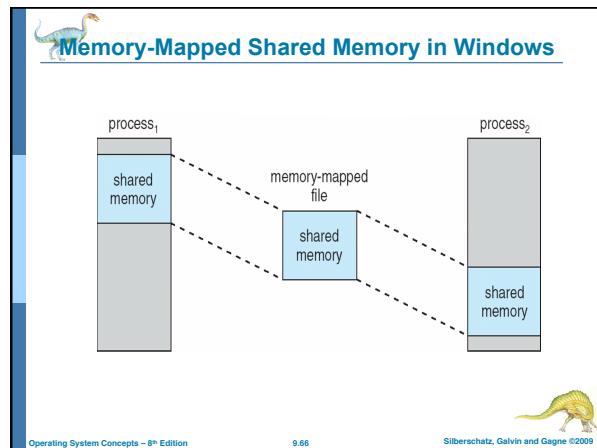
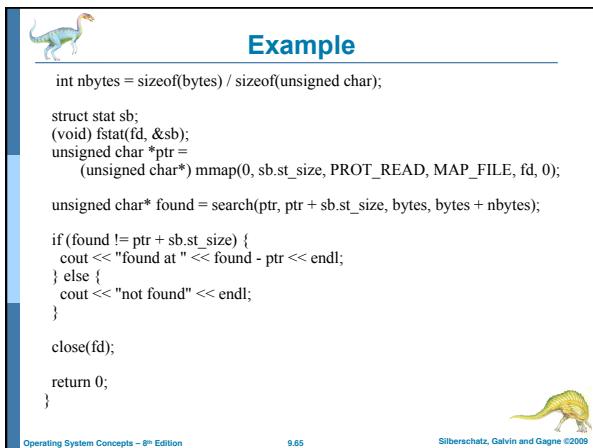
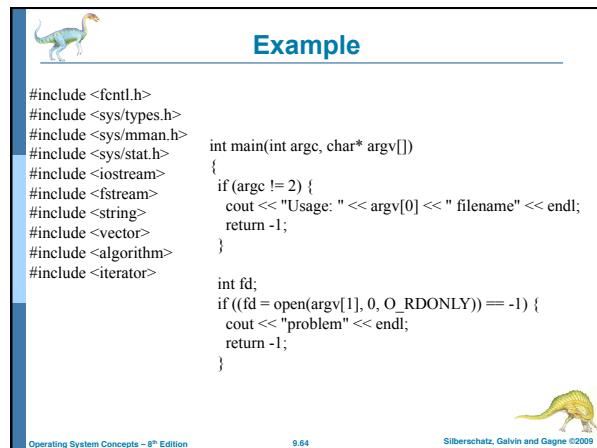
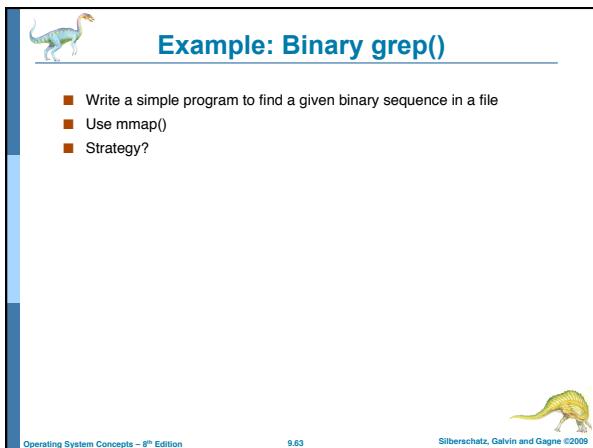
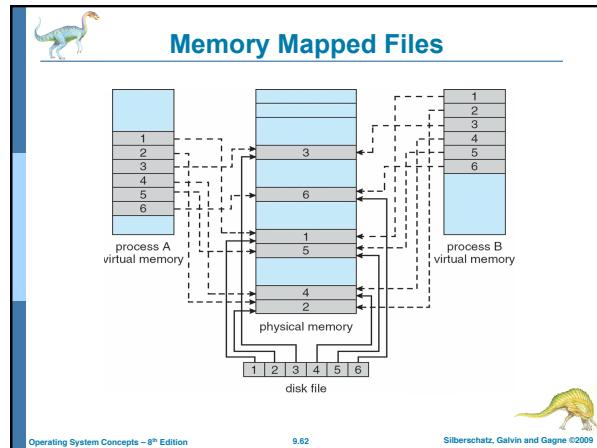
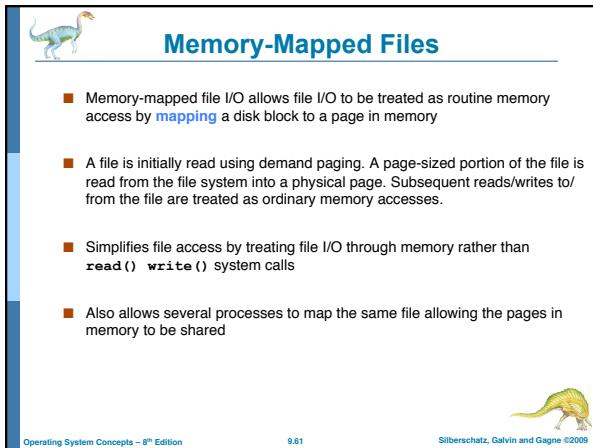


- ### Working-Set Model
- Δ = working-set window = a fixed number of page references
Example: 10,000 instruction
 - WSS_i (working set of Process P_i) = total number of pages referenced in the most recent Δ (varies in time)
 - if Δ too small will not encompass entire locality
 - if Δ too large will encompass several localities
 - if $\Delta = \infty \Rightarrow$ will encompass entire program
 - $D = \sum WSS_i$ = total demand frames
 - if $D > m \Rightarrow$ Thrashing
 - Policy if $D > m$, then suspend one of the processes
- Operating System Concepts – 8th Edition 9.56 Silberschatz, Galvin and Gagne ©2009



- ### 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 interrupt occurs, copy and sets 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
- Operating System Concepts – 8th Edition 9.58 Silberschatz, Galvin and Gagne ©2009





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

Operating System Concepts – 8th Edition 9.67 Silberschatz, Galvin and Gagne ©2009

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

Operating System Concepts – 8th Edition 9.68 Silberschatz, Galvin and Gagne ©2009

Buddy System Allocator

Operating System Concepts – 8th Edition 9.69 Silberschatz, Galvin and Gagne ©2009

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

Operating System Concepts – 8th Edition 9.70 Silberschatz, Galvin and Gagne ©2009

Slab Allocation

Operating System Concepts – 8th Edition 9.71 Silberschatz, Galvin and Gagne ©2009

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

Operating System Concepts – 8th Edition 9.72 Silberschatz, Galvin and Gagne ©2009



Other Issues – Page Size

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

Operating System Concepts – 8th Edition 9.73 Silberschatz, Galvin and Gagne ©2009



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

Operating System Concepts – 8th Edition 9.74 Silberschatz, Galvin and Gagne ©2009



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

Operating System Concepts – 8th Edition 8.75 Silberschatz, Galvin and Gagne ©2009



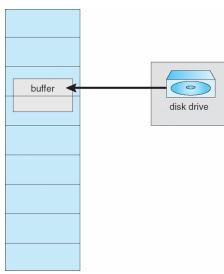
Other Issues – I/O interlock

- 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

Operating System Concepts – 8th Edition 8.76 Silberschatz, Galvin and Gagne ©2009



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



The diagram illustrates the relationship between memory and disk storage. A vertical stack of blue rectangular boxes represents memory frames. At the top of this stack is a smaller red box labeled "buffer". An arrow points from the "buffer" box to a small icon of a hard disk drive, indicating that the buffer is used as an intermediate storage area between the disk and the main memory frames.

Operating System Concepts – 8th Edition 9.77 Silberschatz, Galvin and Gagne ©2009



Operating System Examples

- Windows XP
- Solaris

Operating System Concepts – 8th Edition 9.78 Silberschatz, Galvin and Gagne ©2009



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

Operating System Concepts – 8th Edition

9.79

Silberschatz, Galvin and Gagne ©2009

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 begin 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

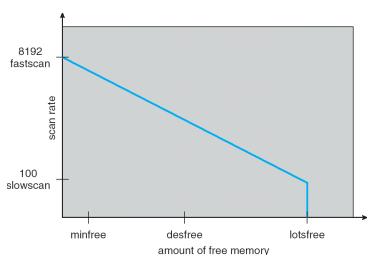
Operating System Concepts – 8th Edition

9.80

Silberschatz, Galvin and Gagne ©2009



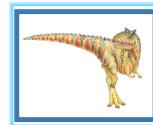
Solaris 2 Page Scanner

Operating System Concepts – 8th Edition

9.81

Silberschatz, Galvin and Gagne ©2009

End of Chapter 9

Operating System Concepts – 8th Edition,

Silberschatz, Galvin and Gagne ©2009