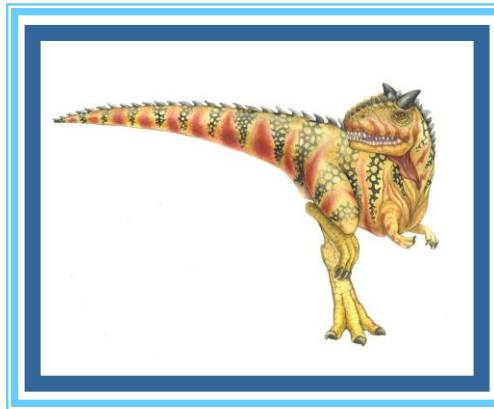


Chapter 10: Virtual Memory





Chapter 10: Virtual Memory

- Background
- Demand Paging
- Copy-on-Write
- Page Replacement
- Allocation of Frames
- Thrashing





Objectives

- Define virtual memory and **describe its benefits**.
- Illustrate how pages are loaded into memory using **demand paging**.
- Apply the FIFO, optimal, and LRU page-replacement algorithms.
- Describe the working **set of a process**, and explain how it is related to **program locality**.
- *Design a virtual memory manager simulation in the C programming language.*





Background

- Code needs to be in **memory to execute**, but entire program rarely used
 - Error code, *unusual routines*, large data structures
- Entire program code **not needed at same time**
- Consider ability to **execute partially-loaded program**
 - Program no longer constrained by **limits of physical memory**
 - Each program takes **less memory while running** -> more programs run at the same time
 - ▶ Increased CPU utilization and throughput with no increase in **response time or turnaround time**
 - Less **I/O needed to load or swap programs into memory** -> each user program runs faster





Virtual memory

- **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** (`vfork()`)
 - More programs running concurrently (Degree of multiprogramming)





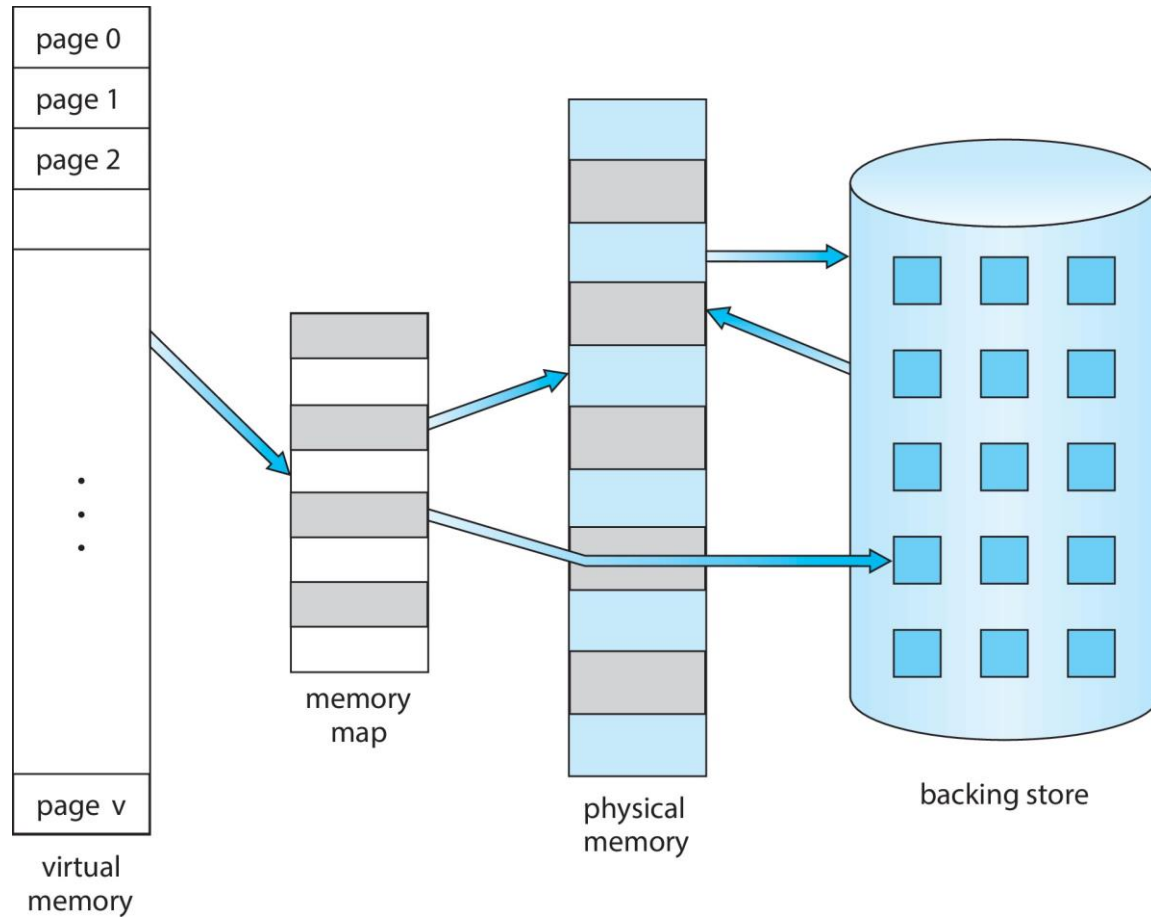
Virtual memory (Cont.)

- **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**





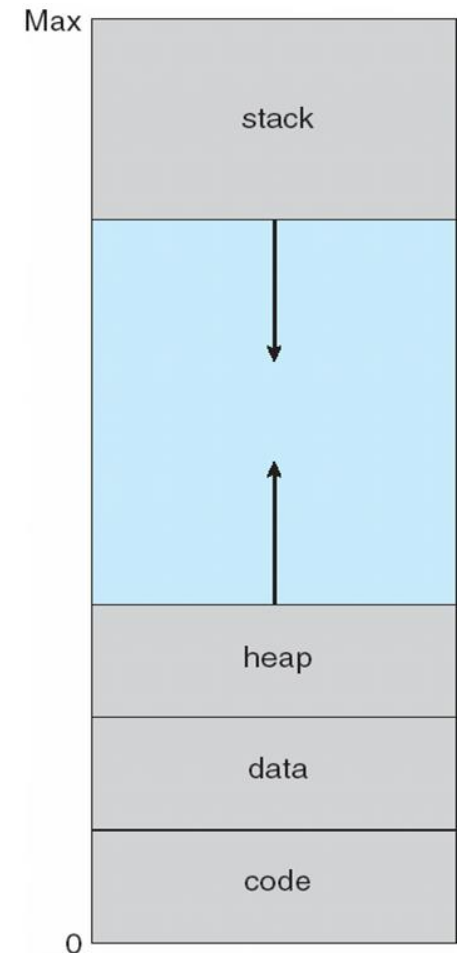
Virtual Memory That is Larger Than Physical Memory

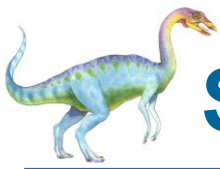




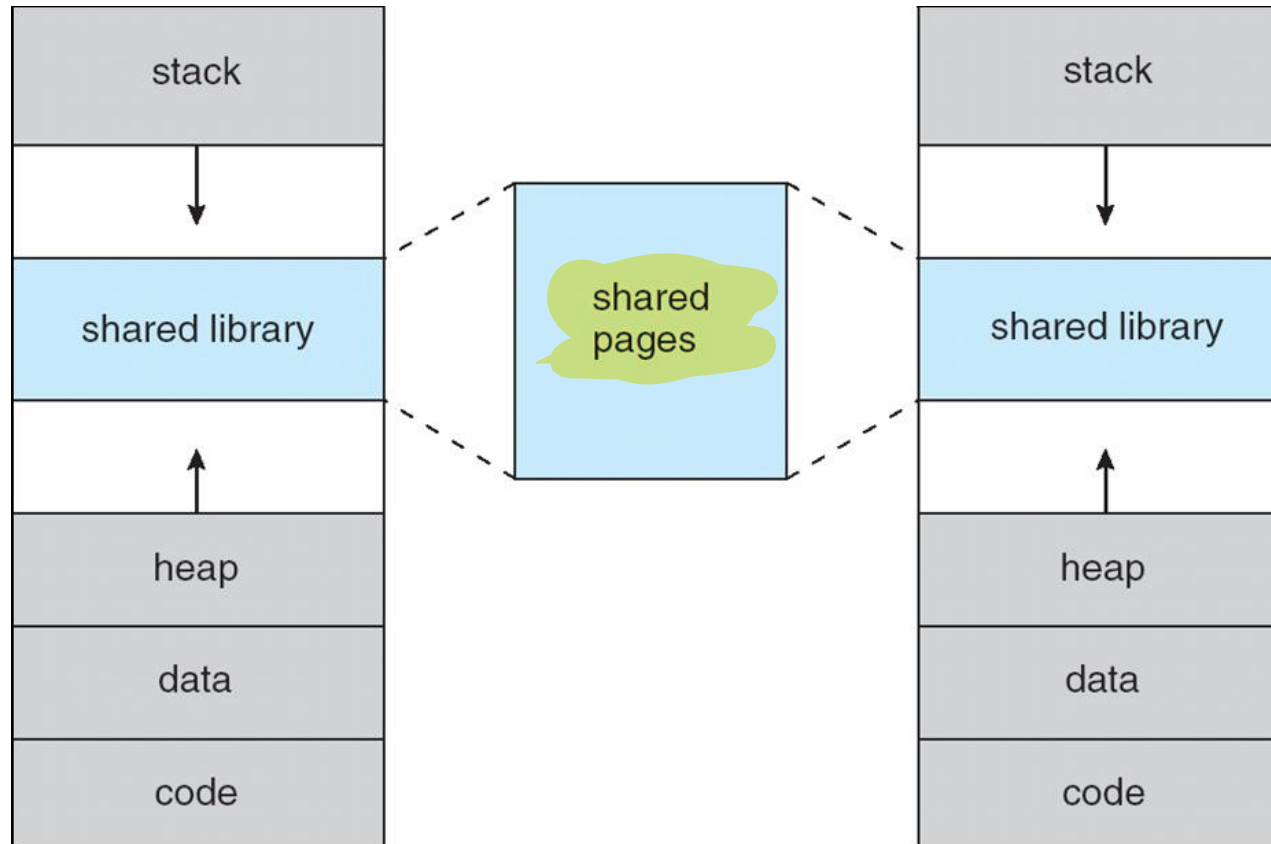
Virtual-address Space

- Usually design *logical address space for stack* to start at *Max logical address* and grow “down” while heap grows “up”
 - Maximizes address space use
 - Unused address space between the two is hole
 - ▶ No physical memory needed until heap or stack grows to a given new page
- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc.
- System libraries shared via mapping into virtual address space
 - Although each process considers the libraries to be part of its virtual address space, the actual pages where the libraries reside in physical memory are shared by all the processes





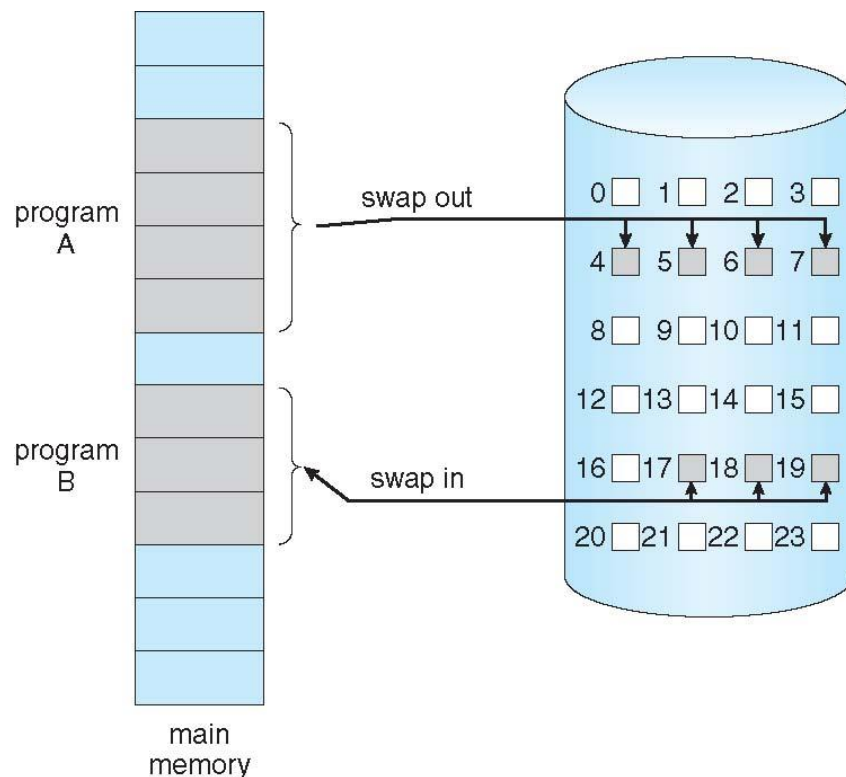
Shared Library Using Virtual Memory





Demand Paging

- 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





Valid-Invalid Bit

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

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

page table

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





Page Table When Some Pages Are Not in Main Memory

0	A
1	B
2	C
3	D
4	E
5	F
6	G
7	H

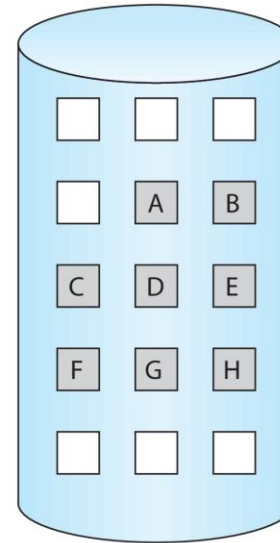
logical memory

valid-invalid bit		
frame		bit
0	4	v
1		i
2	6	v
3		i
4		i
5	9	v
6		i
7		i

page table

0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

physical memory

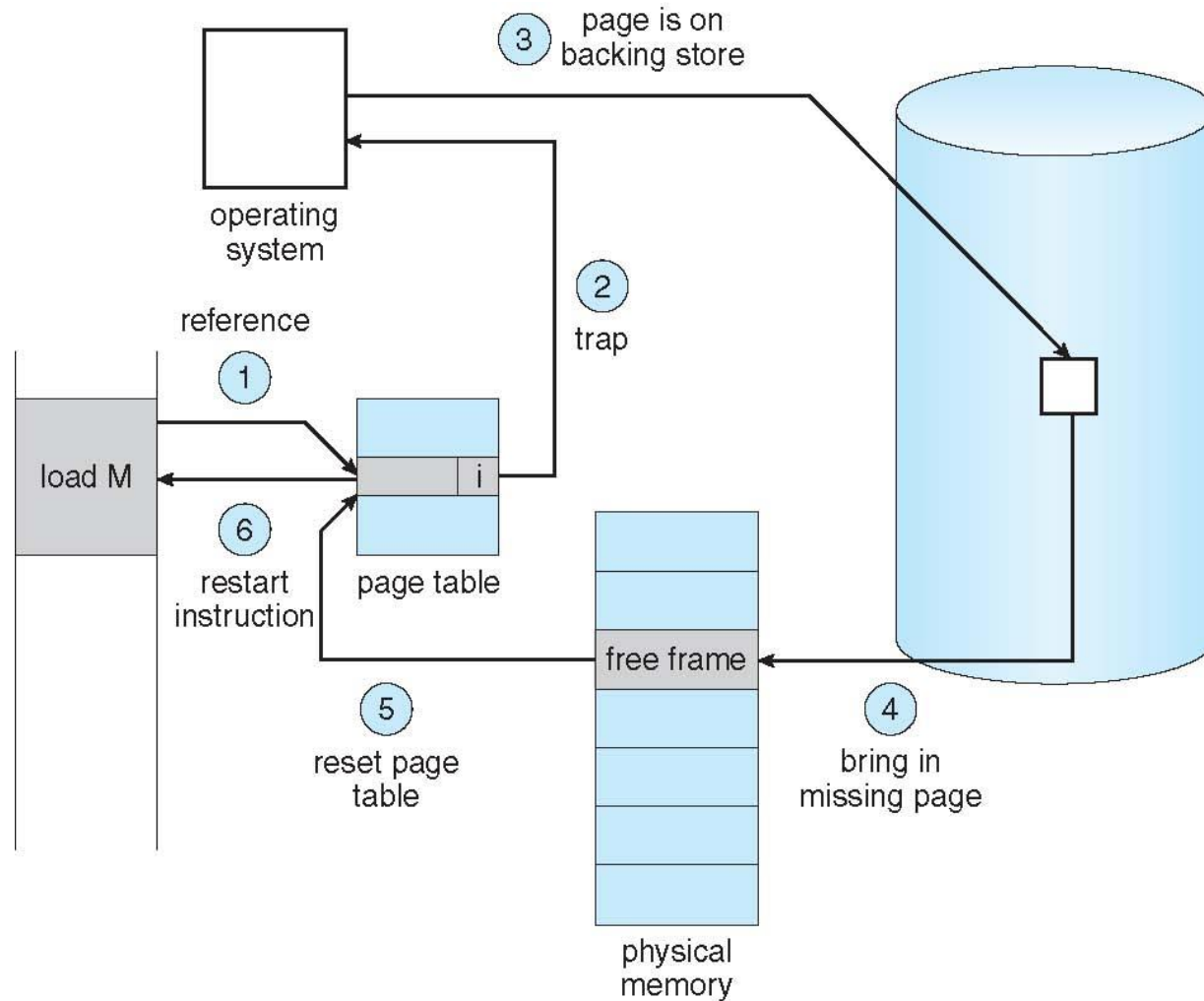


backing store





Steps in Handling a Page Fault (Cont.)





Steps in Handling Page Fault

1. If there is a reference to a page, first reference to that page **will trap to operating system (Why??)**
 - **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





Aspects of Demand Paging

- Extreme case – **start process with *no* pages in memory**
 - OS sets instruction pointer to first instruction of process, non-memory-resident -> page fault
 - And for every other process pages on first access
 - **Pure demand paging**
- Actually, a given *instruction could access multiple pages* -> multiple page faults
 - Consider fetch and decode of instruction which adds 2 numbers from memory and stores result back to memory
 - Pain decreased because of **locality of reference** temporal and spatial locality
- 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.



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





Stages in Demand Paging – Worse Case

- Demand paging can significantly affect the performance of a computer system. (consider the given situation!!)
- 1. Trap to the operating system
- 2. Save the user registers and process state (Why??)
- 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:
 - a) Wait in a queue for this device until the read request is serviced
 - b) Wait for the device seek and/or latency time
 - c) 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)

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





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.
This is a slowdown by a factor of 40!!
- If want performance degradation < 10 percent
 - $220 > 200 + 7,999,800 \times p$
 $20 > 7,999,800 \times p$
 - $p < .0000025$
 - < one page fault in every 400,000 memory accesses
- What is the take away??





Demand Paging Optimizations

- Swap space I/O faster than file system I/O even if on the same device
 - Swap allocated in larger chunks, less management needed than file system
- Copy entire process image to swap space at process load time
 - Then page in and out of swap space
 - Used in older BSD Unix
- Demand page in from program binary on disk, but discard rather than paging out when freeing frame
 - Used in Solaris and current BSD
 - Still need to write to swap space
 - ▶ Pages not associated with a file (like stack and heap) – **anonymous memory**
 - ▶ Pages modified in memory but not yet written back to the file system
- Mobile systems
 - Typically don't support swapping
 - Instead, demand page from file system and reclaim read-only pages (such as code)





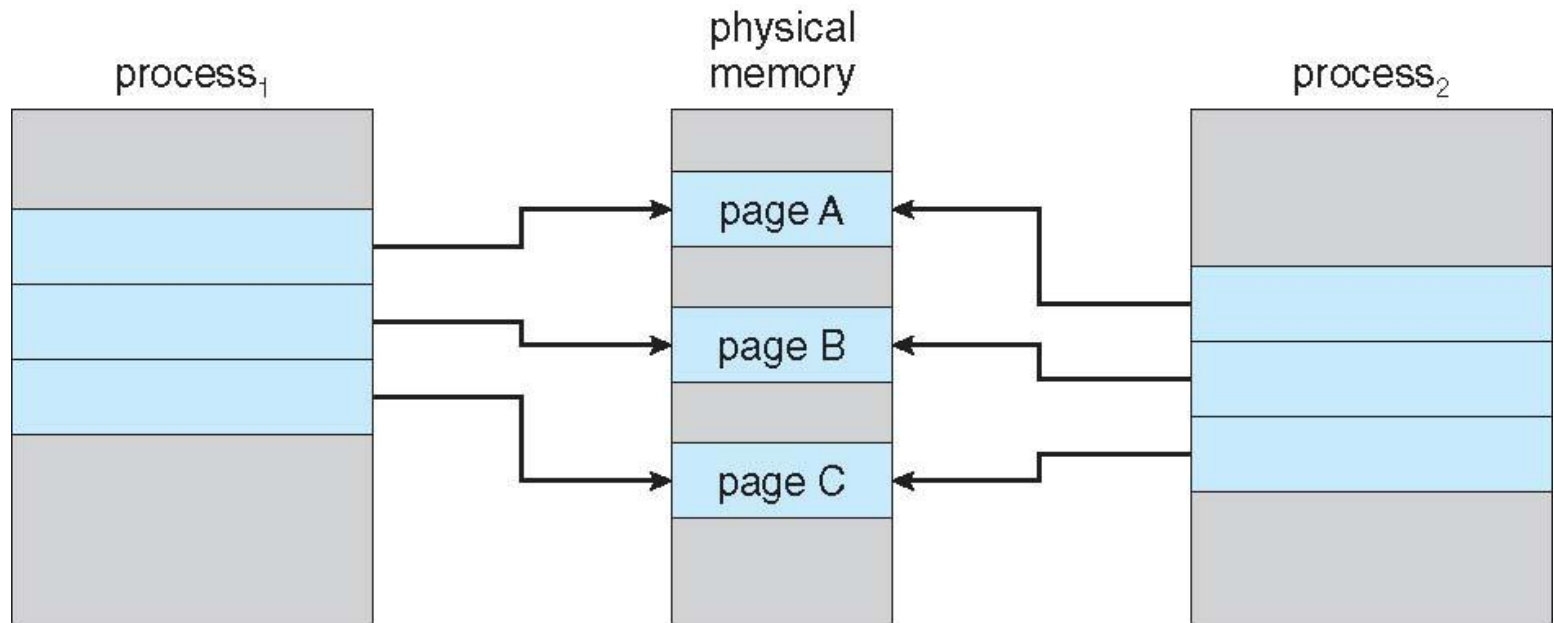
Copy-on-Write

- **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
- In general, free pages are allocated from a **pool** of **zero-fill-on-demand** pages
 - Pool should always have free frames for fast demand page execution
 - ▶ Don't want to have to free a frame as well as other processing on page fault
 - Why zero-out a page before allocating it?



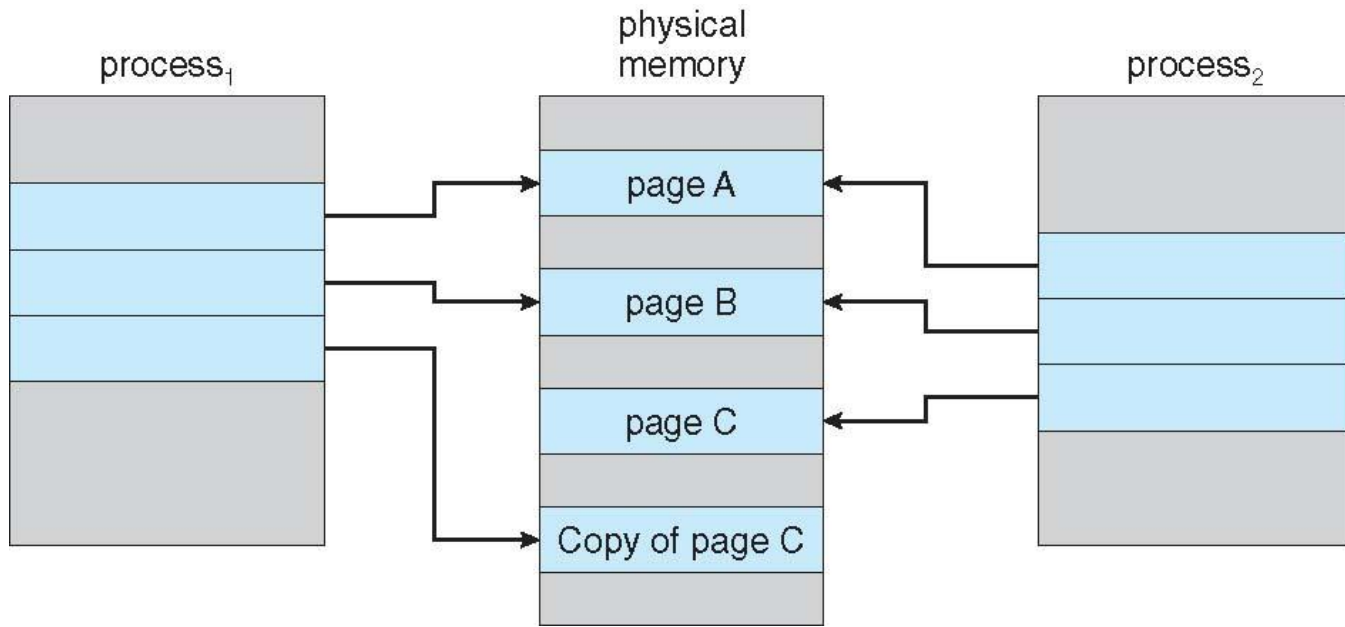


Before Process 1 Modifies Page C





After Process 1 Modifies Page C





fork() and vfork()

- fork() system call creates a child process that is a duplicate of its parent.
 - Traditionally, fork() worked by creating a copy of the parent's address space for the child
 - many child processes invoke the exec() system call immediately after creation, the copying of the parent's address space may be unnecessary.
 - fork() with copy-on-write technique can be used.
- virtual memory fork i.e., vfork()
 - With vfork(), the parent process is suspended, and the child process uses the address space of the parent.
 - Because vfork() does not use copy-on-write, if the child process changes any pages of the parent's address space, the altered pages will be visible to the parent once it resumes.
 - vfork() must be used with caution to ensure that the child process does not modify the address space of the parent.





What Happens if There is no Free Frame?

- Used up by process pages
- Also in demand from the kernel, I/O buffers, etc
- Deciding **how much memory to allocate** to I/O and how much to program pages is a significant challenge?
- **Page replacement** – find some page in memory, but not really in use, page it out
 - Algorithm – terminate process? swap out process? 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





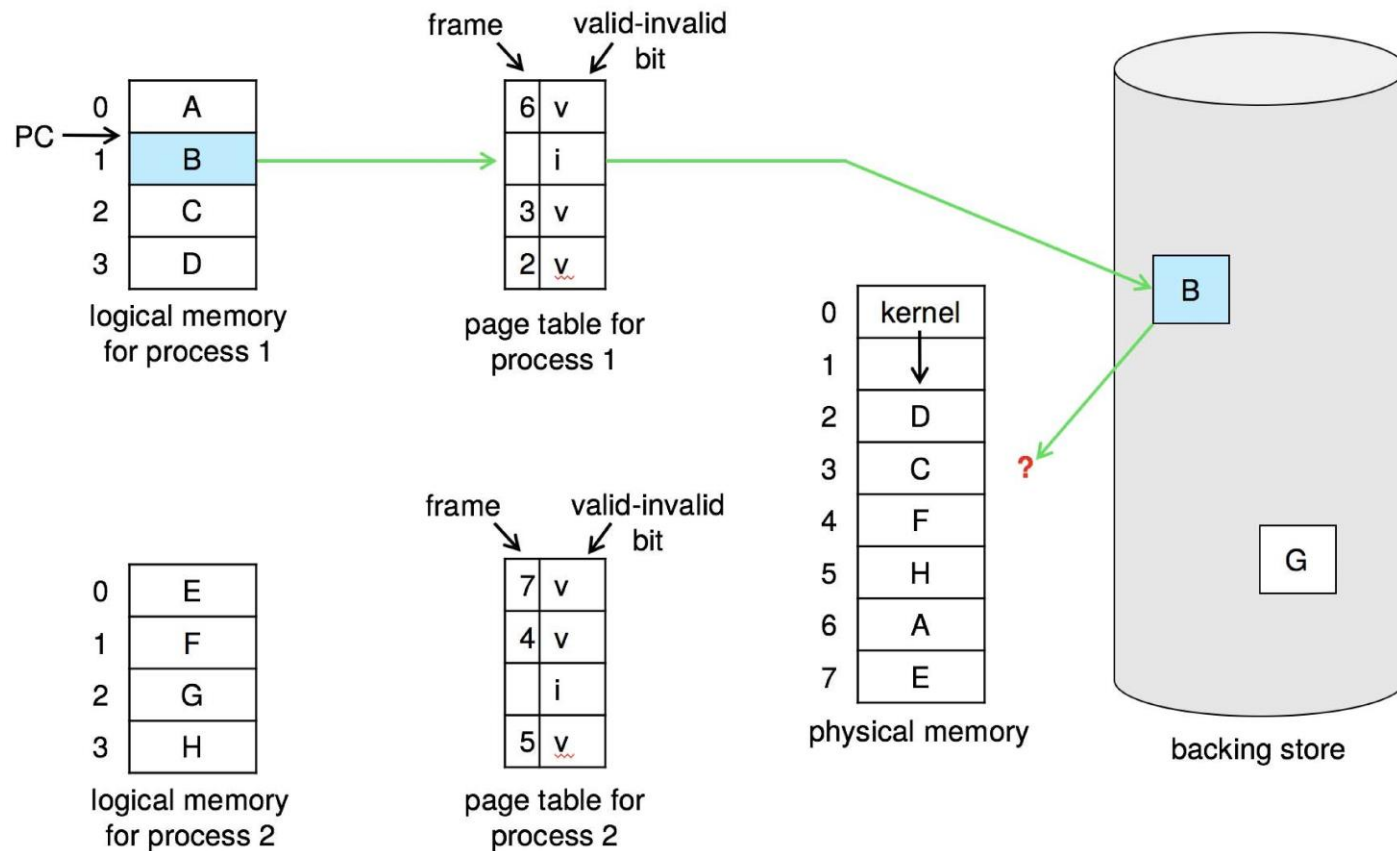
Page Replacement

- 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**





Need For Page Replacement





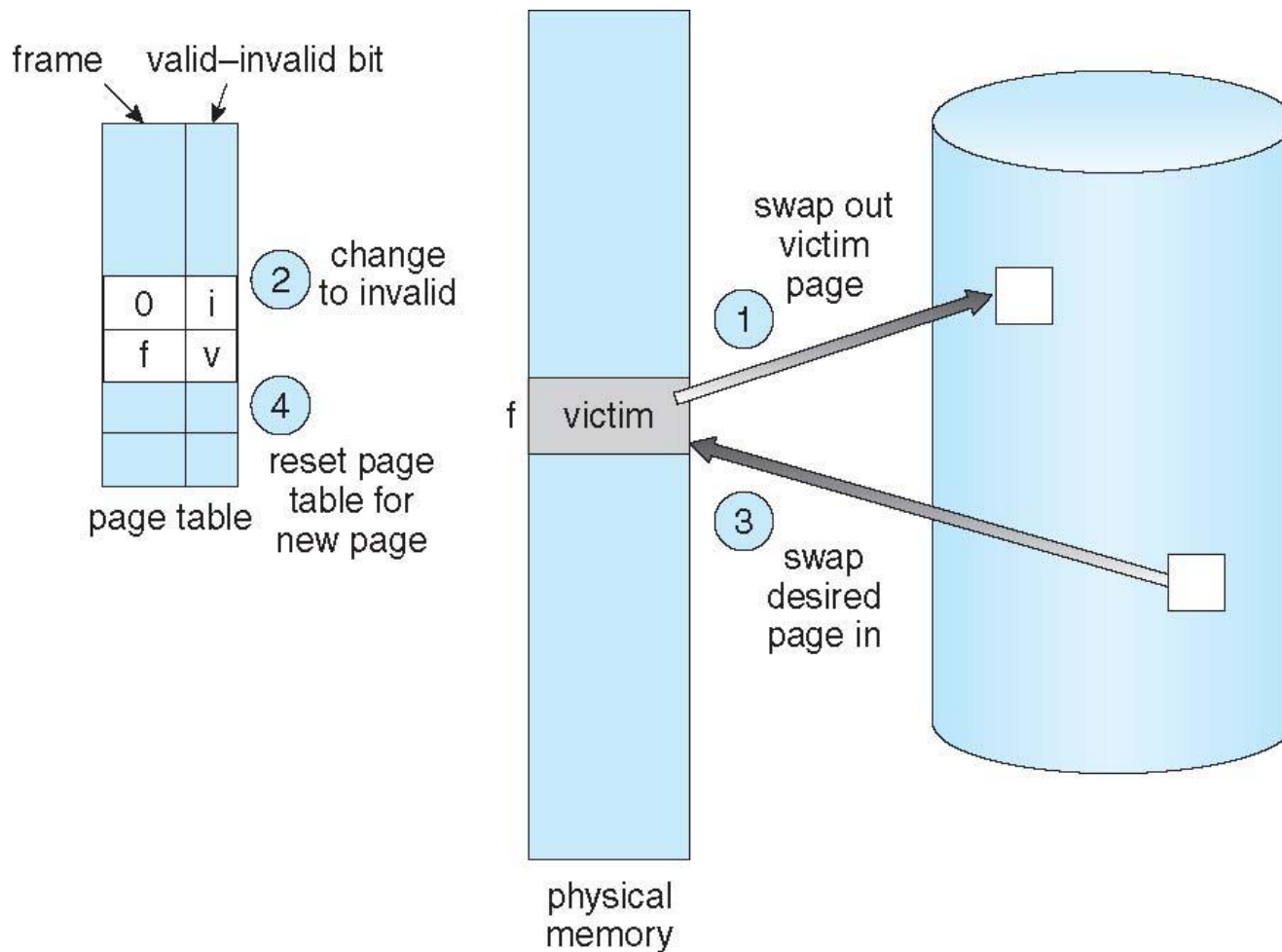
Basic Page Replacement

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





Page Replacement





Page and Frame Replacement Algorithms

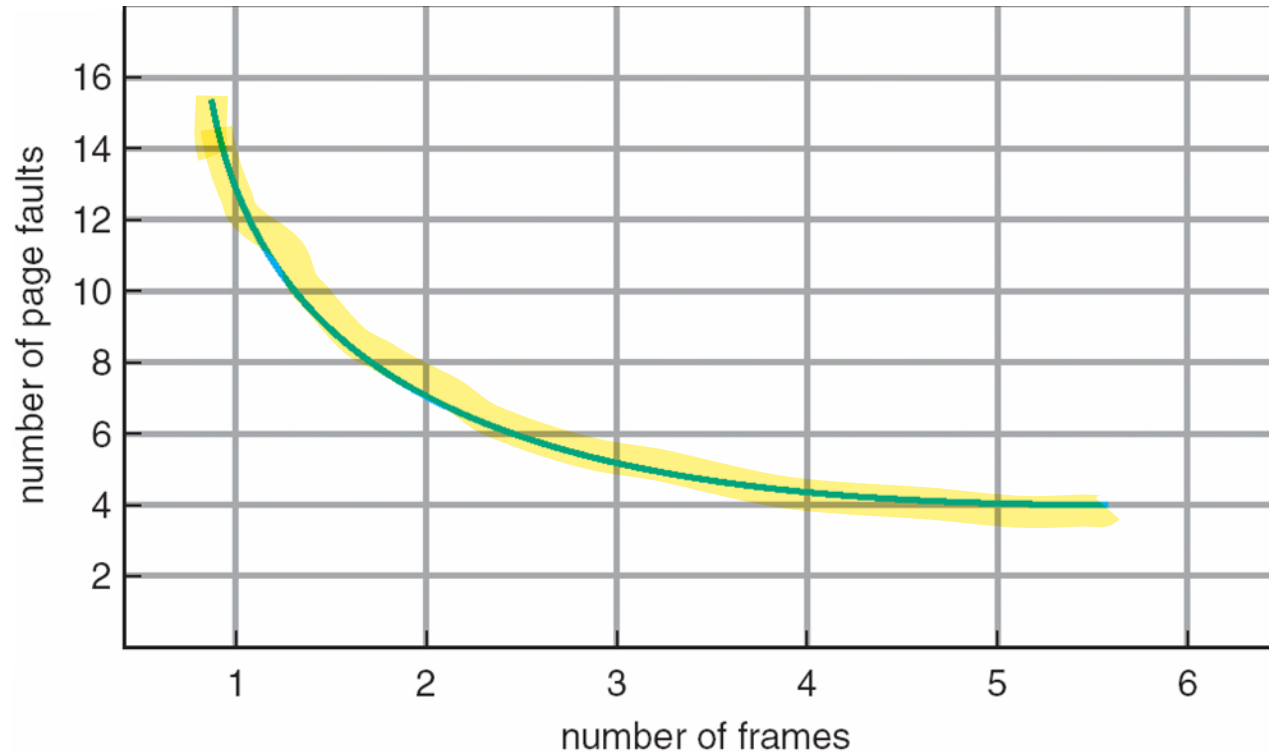
- **Frame-allocation algorithm** determines
 - How many frames to give each process
 - Which frames to replace
- **Page-replacement algorithm**
 - Want **lowest page-fault rate** on **both first access and re-access**
- Evaluate algorithm by running it on a particular string of memory references (reference 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
- In all our examples, the **reference string** of referenced page numbers is

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





Graph of Page Faults Versus the Number of Frames





First-In-First-Out (FIFO) Algorithm

- Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,2,1,2,0,1,7,0,1**
- 3 frames (**3 pages can be in memory at a time per process**)

reference string

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

7	7	7	2																
	0	0	0																
		1	1																

page frames

15 page faults

- How to track ages of pages?
 - Just use a FIFO queue

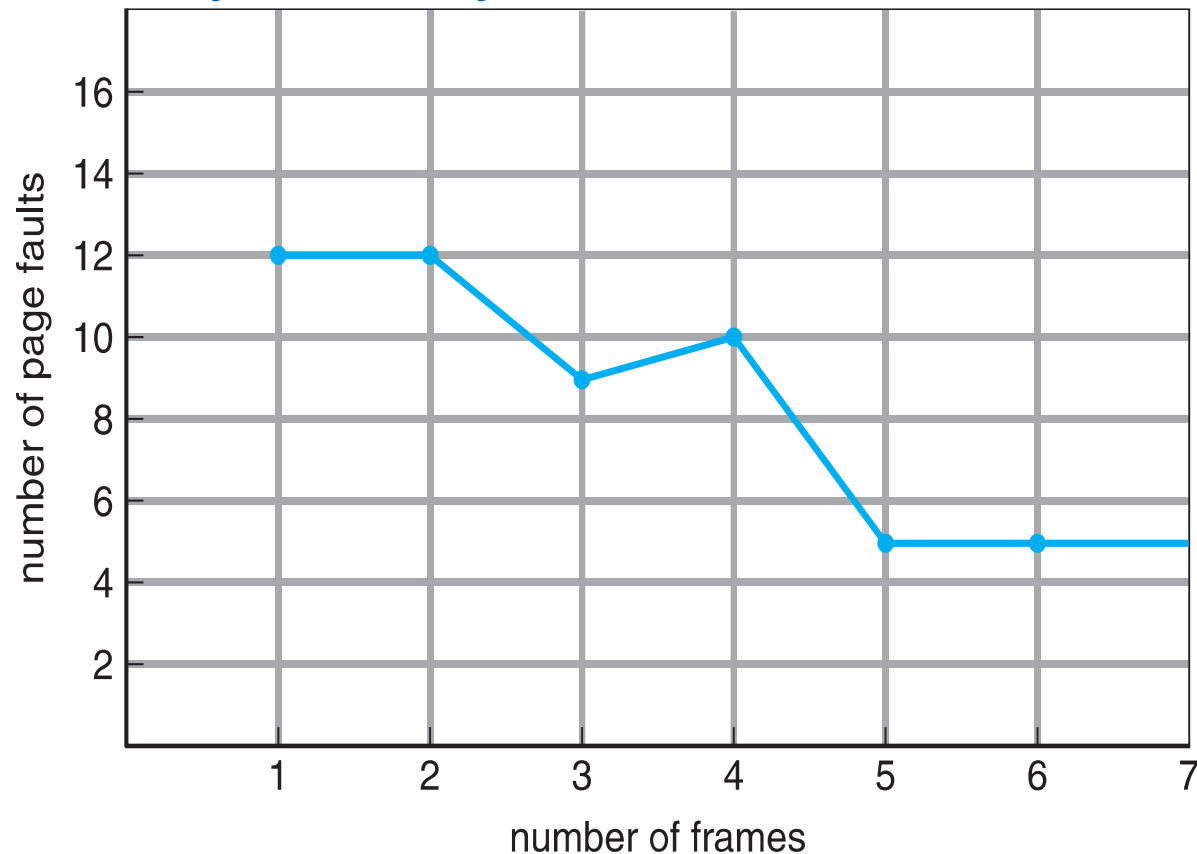




FIFO Illustrating Belady's Anomaly

- # page fault vary by reference string: consider 1,2,3,4,1,2,5,1,2,3,4,5
- Adding more frames can cause more page faults!

► Belady's Anomaly





Optimal Algorithm

- Replace page that will not be used for **longest period of time**
 - 9 is optimal for the example
- How do you know this?
 - Can't read the future
- Used for measuring **how well your algorithm performs**

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		0		0		0					0		
		1	1		3		3		3		1					1		

page frames





Least Recently Used (LRU) Algorithm

- Use past knowledge rather than future
- Replace page that has not been used in the most amount of time
- Associate time of last use with each page

reference string

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

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

page frames

- 12 faults – better than FIFO but worse than OPT
- 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**
 - **But each update more expensive**
 - **No search for replacement**





LRU Algorithm (Cont.)

- LRU and OPT are cases of **stack algorithms** that don't have **Belady's Anomaly (Justification??)**

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

↑ ↑
a b





LRU Algorithm (Cont.)

- LRU and OPT are cases of **stack algorithms** that don't have **Belady's Anomaly (Justification??)**
- A **stack algorithm** is an algorithm for which it can be shown that 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.

reference string

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

2
1
0
7
4

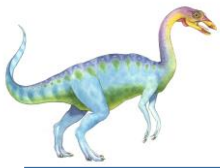
stack
before
a

7
2
1
0
4

stack
after
b

↑ ↑
a b





LRU Approximation Algorithms

- Unfortunately full implementation of LRU requires hardware support, and **few systems provide this**.
- However many systems offer some degree of HW support, enough to approximate LRU fairly well.
 - many systems provide a **Reference bit** for every entry in a page table
 - ▶ initially = 0
 - ▶ When page is referenced bit set to 1
 - ▶ One bit of precision is enough to distinguish pages that have been accessed since the last clear
 - does not provide any finer grain of detail
 - » We do not know the order, however





Additional-Reference-Bits Algorithm

- Finer grain is possible by storing the most recent **8 reference bits (unsigned int)** for each page in the page table entry.
 - At periodic intervals, the OS takes over, and **right-shifts** each of the reference bytes by **one bit**.
 - The high-order (leftmost) bit is then **filled in with the current value of the reference bit**, and the **reference bits are cleared**.
 - At any given time, the page with the **smallest value for the reference byte** is the LRU page.





LRU Approximation Algorithms (cont.)

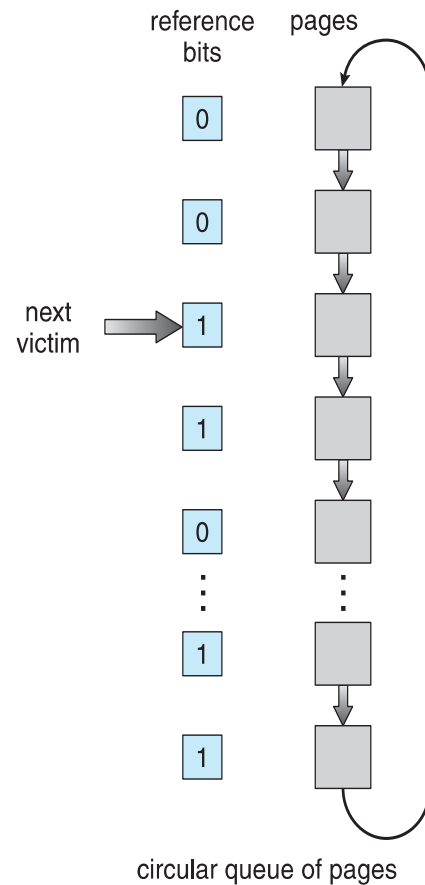
■ Second-chance algorithm

- 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**
- What if all reference bits in the table are set??
 - ▶ **Good as FIFO?**

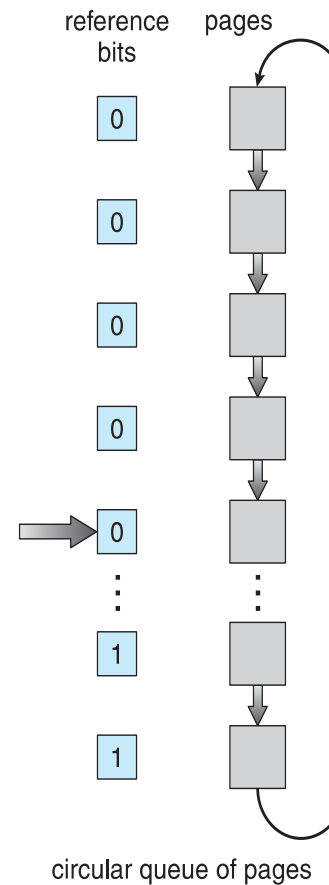




Second-chance Algorithm



(a)



(b)





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 nor modified – **best page to replace (Why??)**
 - (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**





Counting Algorithms

- Keep a counter of the number of references that have been made to each page
 - Not common
- **Least Frequently Used (LFU) Algorithm:**
 - Replaces page with smallest count
- **Most Frequently Used (MFU) Algorithm:**
 - 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

- Maintain a certain minimum number of **free frames at all times**.
 - When a page-fault occurs, go ahead and allocate one of the free frames from the free list first, to **up the requesting process ASAP**.
 - select a victim page to **write to disk and free up a frame as a second step**.
- Possibly, keep list of modified pages
 - When backing store otherwise idle, write pages there and set to non-dirty
- Possibly, keep free **frame contents intact** and note what is in them
 - If referenced again before reused, no need to load contents again from disk
 - Generally useful to reduce penalty if wrong victim frame selected





Allocation of Frames

- Each process needs **minimum** number of frames
- 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*
- **Maximum** of course is total frames in the system
- Two major allocation schemes
 - fixed allocation
 - priority allocation
- Many variations





Fixed Allocation

- Equal allocation – For example, if there are 100 frames (after allocating frames for the OS) and 5 processes, give each process 20 frames
 - Keep some as **free frame buffer pool**
- **Proportional allocation** – Allocate according to the size of process
 - Dynamic as degree of multiprogramming, process sizes change

– s_i = size of process p_i

– $S = \sum s_i$

– m = total number of frames

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

$$m = 64$$

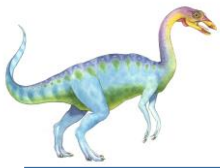
$$s_1 = 10$$

$$s_2 = 127$$

$$a_1 = \frac{10}{137} \cdot 62 \gg 4$$

$$a_2 = \frac{127}{137} \cdot 62 \gg 57$$





Global vs. Local Allocation

- **Global replacement** – process selects a replacement frame from the set of all frames; **one process can take a frame from another**
 - But then process execution time can vary greatly
 - But **greater throughput** so more common
- **Local replacement** – each process selects from only its own set of allocated frames
 - More consistent per-process performance
 - But possibly underutilized memory





Practice Problem

- Consider the following page reference string:
 - 7, 2, 3, 1, 2, 5, 3, 4, 6, 7, 7, 1, 0, 5, 4, 6, 2, 3, 0, 1.

Assuming demand paging with three frames, how many page faults would occur for the following replacement algorithms?

- LRU replacement
- FIFO replacement
- Optimal replacement
- Second Chance Algorithm





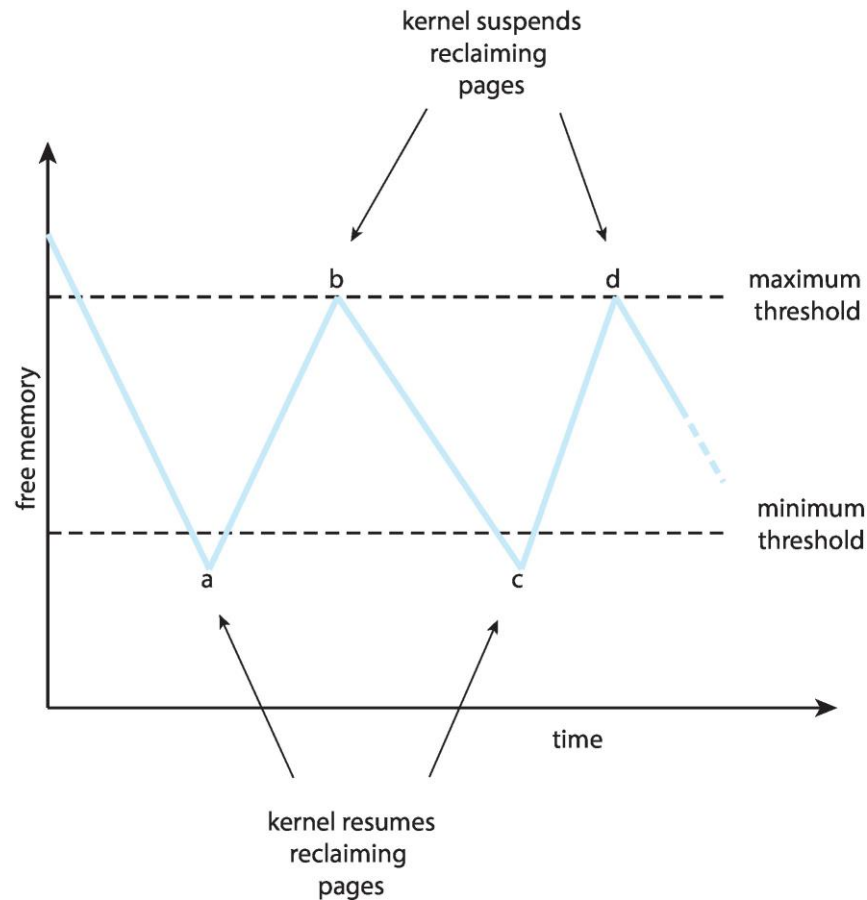
Reclaiming Pages

- A strategy to implement **global page-replacement policy**
- All memory requests are satisfied from the **free-frame list**, rather than waiting for the list to drop to zero before we begin selecting pages for replacement,
- Page replacement is triggered when the **list falls below a certain threshold**.
- This strategy attempts to ensure there is **always sufficient free memory to satisfy new requests**.





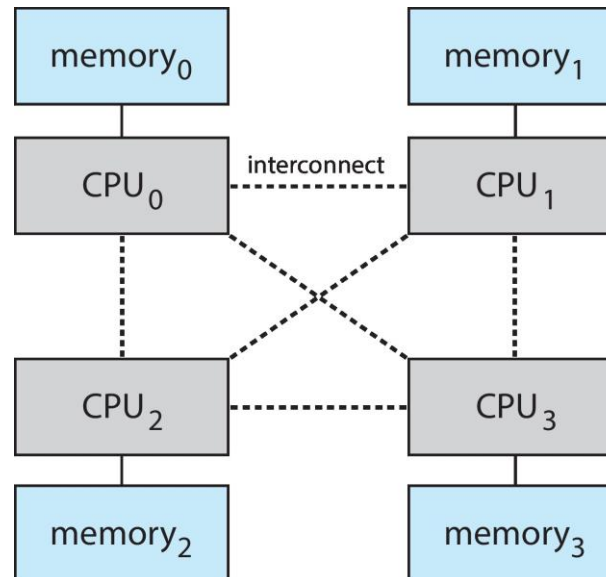
Reclaiming Pages Example





Non-Uniform Memory Access

- So far, we assumed that all memory accessed equally
- Many systems are **NUMA** – speed of access to memory varies
 - Consider system boards containing CPUs and memory, interconnected over a system bus
- NUMA multiprocessing architecture





Thrashing

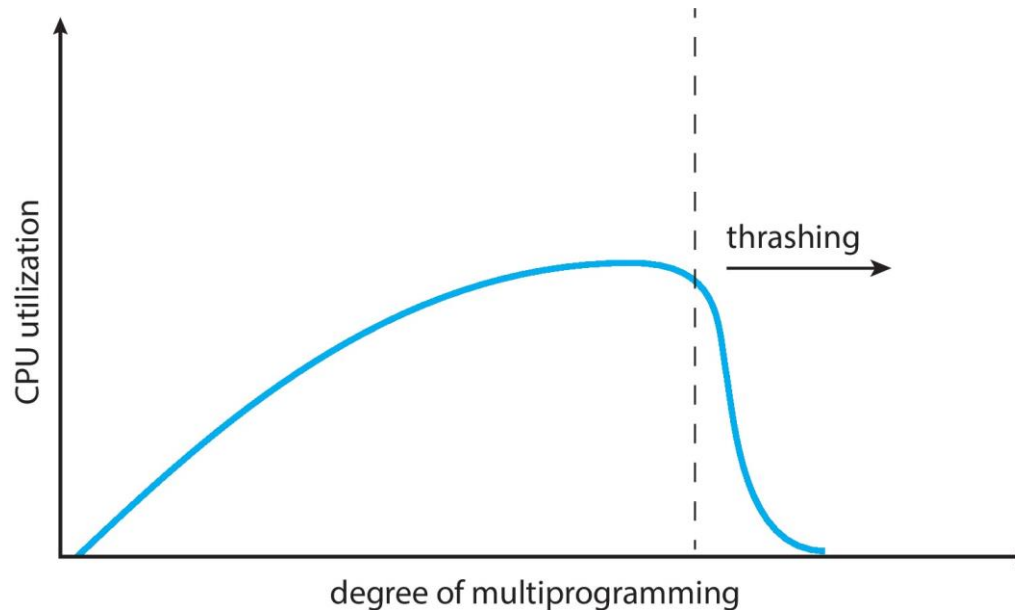
- If a process does not have “enough” pages, the page-fault rate is very high
 - Page fault to get page
 - Replace existing frame
 - But quickly need replaced frame back
 - This leads to:
 - ▶ Low CPU utilization
 - ▶ Operating system thinking that it needs to increase the degree of multiprogramming
 - ▶ Another process added to the system





Thrashing (Cont.)

- **Thrashing.** A process is busy swapping pages in and out
- A process is thrashing if it is **spending more time paging than executing**





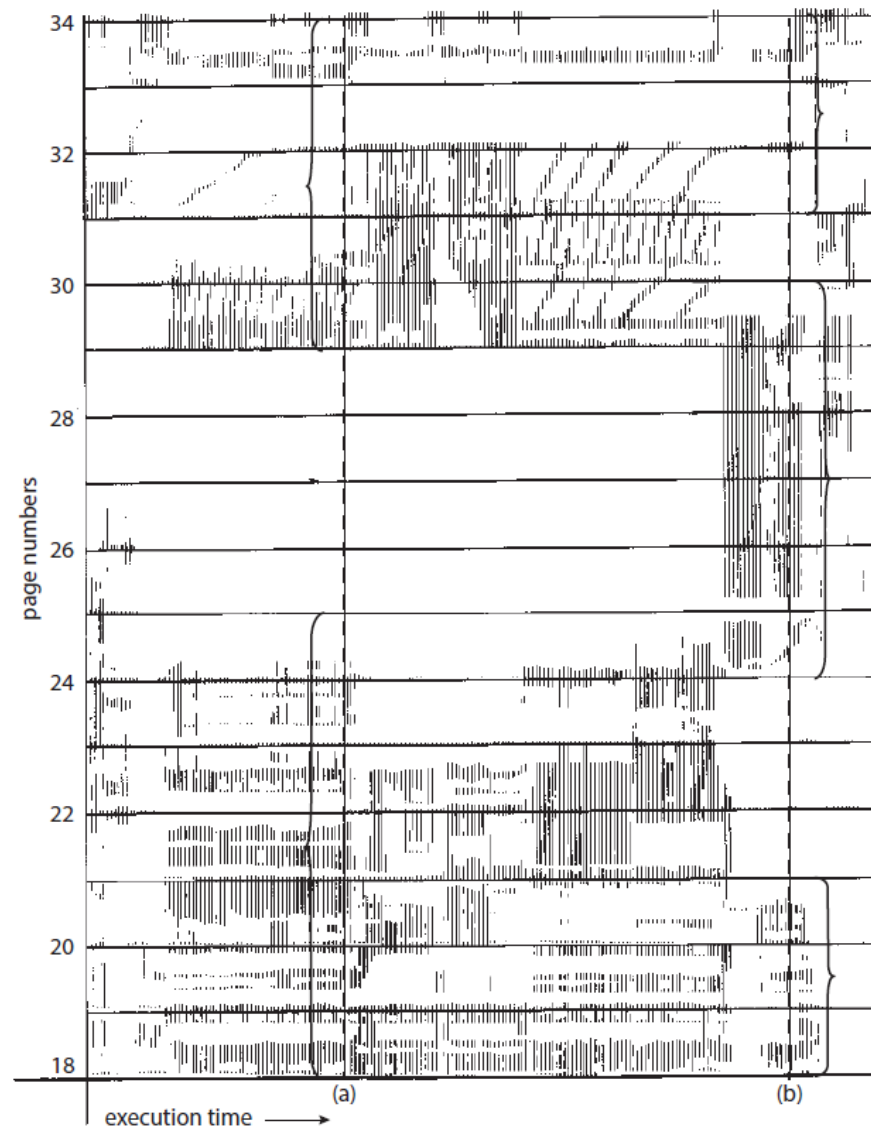
Demand Paging and Thrashing

- To prevent thrashing, we must provide a process with as many frames as it needs. But how do we know how many frames it “needs”?
 - Looking at how many frames a process is actually using.
 - This approach defines the **locality model** of process execution.
- Locality model states that, as a process executes, it moves from locality to locality.
 - Locality is a set of pages that are actively used together.
 - A running program is generally composed of several different localities, which may overlap.
- Why does thrashing occur?
$$\Sigma \text{ size of locality} > \text{total memory size}$$
- Limit effects by using local or priority page replacement





Locality In A Memory-Reference Pattern



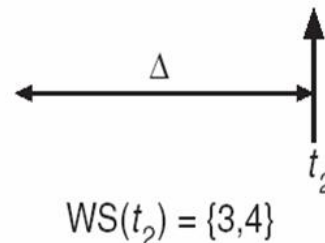
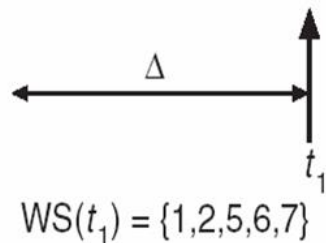


Working-Set Model

- $\Delta \equiv$ working-set window \equiv a **fixed number of page references**
- WSS_i (**working-set size** 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
 - Approximation of locality

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 ...



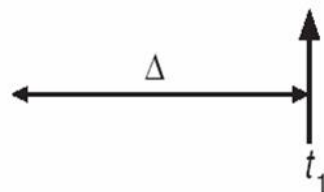


Working-Set Model (Cont.)

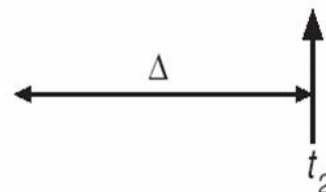
- if $D > m \Rightarrow$ Thrashing
 - m : total number of available frames
- Policy if $D > m$, then **suspend or swap out one of the processes**

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 4 1 3 2 3 4 4 4 3 4 4 4 ...



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



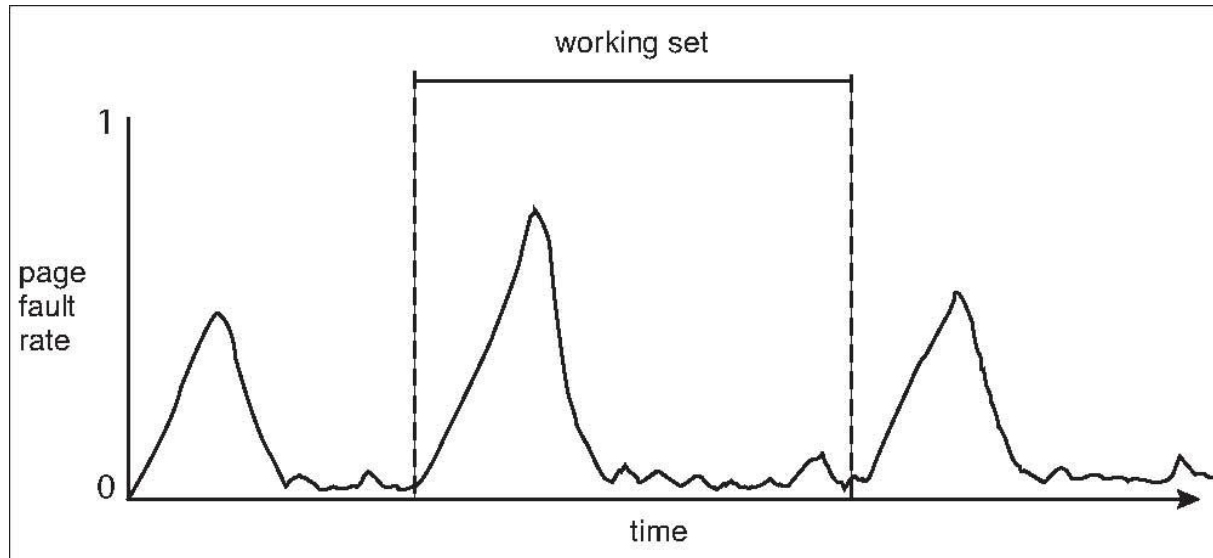
$$WS(t_2) = \{3, 4\}$$





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





Practice Problem

- Assume we have a demand-paged memory. The page table is held in registers. It takes 8 milliseconds to service a page fault if an empty page is available or the replaced page is not modified, and 20 milliseconds if the replaced page is modified. Memory access time is 100 nanoseconds. Assume that the page to be replaced is modified 70 percent of the time. What is the **maximum acceptable page-fault rate** for an effective access time of no more than 200 nanoseconds?





Practice Problem

- Consider a demand-paging system with the following time-measured utilizations: CPU utilization 20%; Paging disk 97.7%; Other I/O devices 5%. Which of the following will (probably) **improve CPU utilization**?
 - Installing a faster CPU.
 - Installing a bigger paging disk.
 - Increasing the degree of multiprogramming.
 - Decreasing the degree of multiprogramming.
 - Installing more main memory.
 - Installing a faster hard disk or multiple controllers with multiple hard disks.



End of Chapter 10

