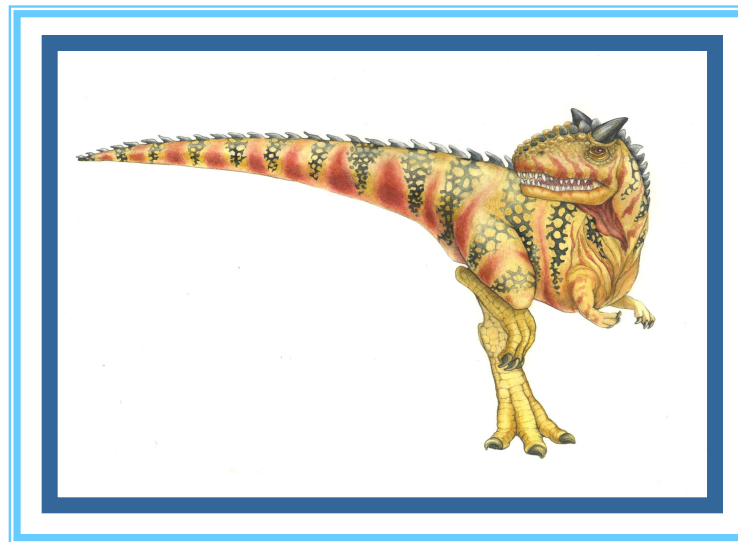


# Chapter 9: Virtual Memory

---



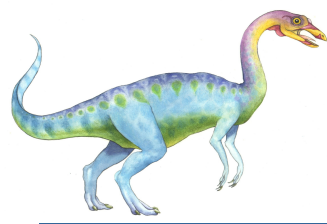


# Background

---

- Program no longer constrained by limits of physical memory
- A program could be larger than physical memory.
- The Logical address space is no longer constrained by the physical address space or the actual physical memory
- Logical address space could be larger than physical memory
- Logical address space could be smaller than the physical address space





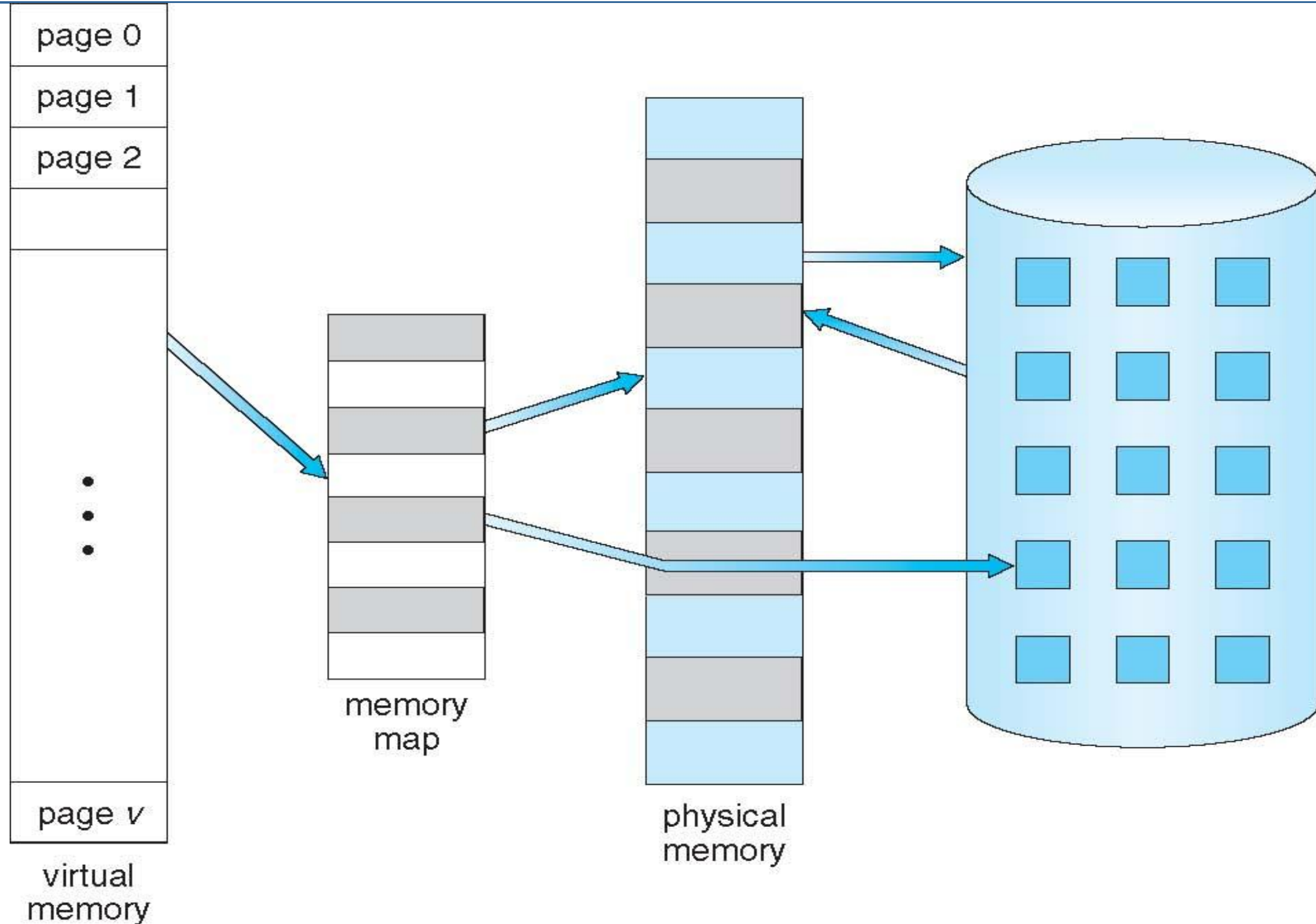
# Background

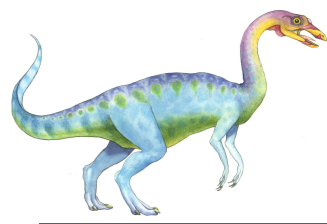
- **Virtual memory** – separation of user logical memory from physical memory
  - Only part of the program needs to be in memory for execution
  - Time to start a large program gets reduced as only a few pages needed in main memory to start the program
  - Allows address spaces to be shared by several processes
  - Allows for more efficient process creation
  - More programs running concurrently
  - Less I/O needed to load or swap processes
- Virtual memory can be implemented via:
  - Demand paging
  - Demand segmentation



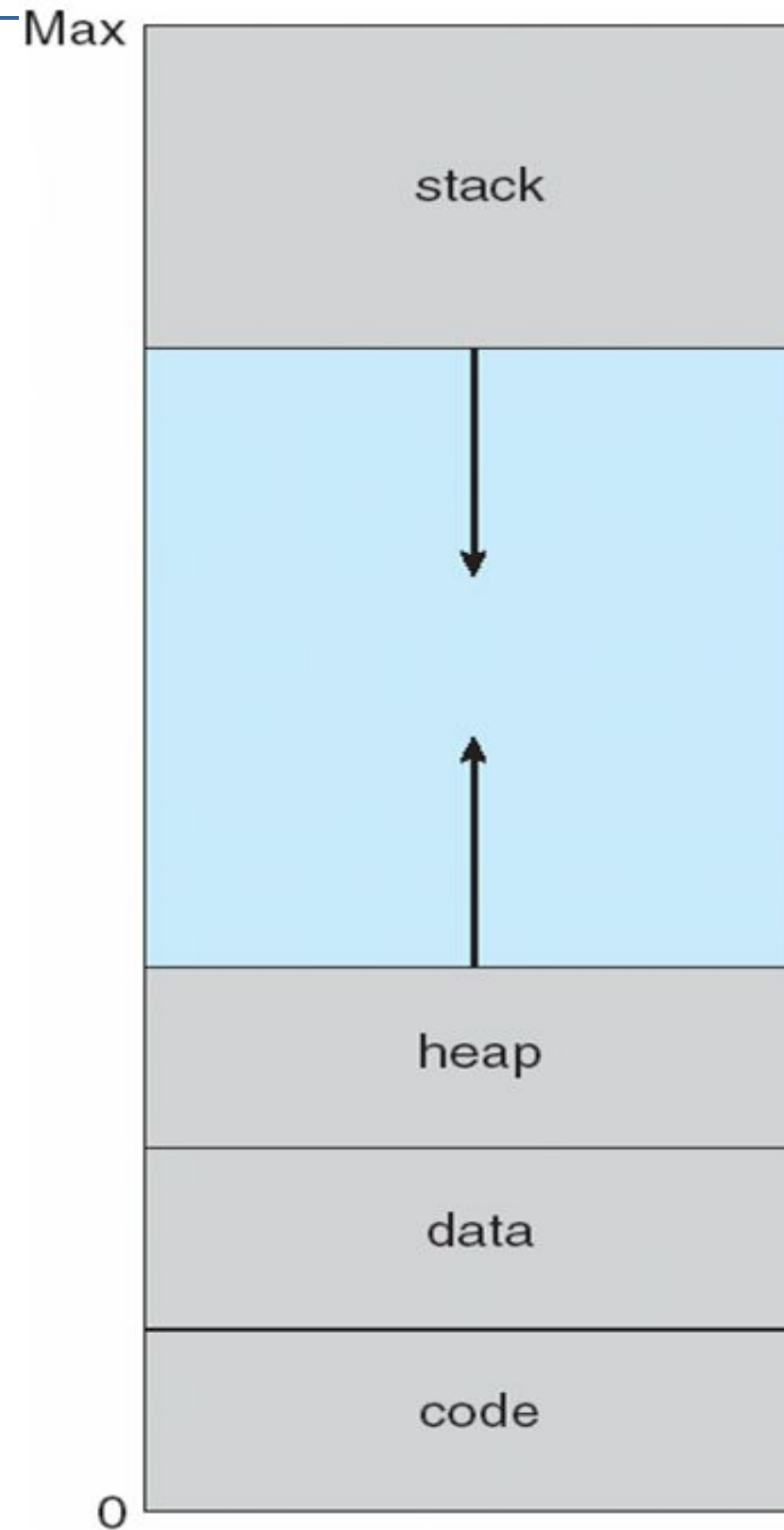


# Virtual Memory That is Larger Than Physical Memory





# Virtual-address Space





# Virtual Address Space

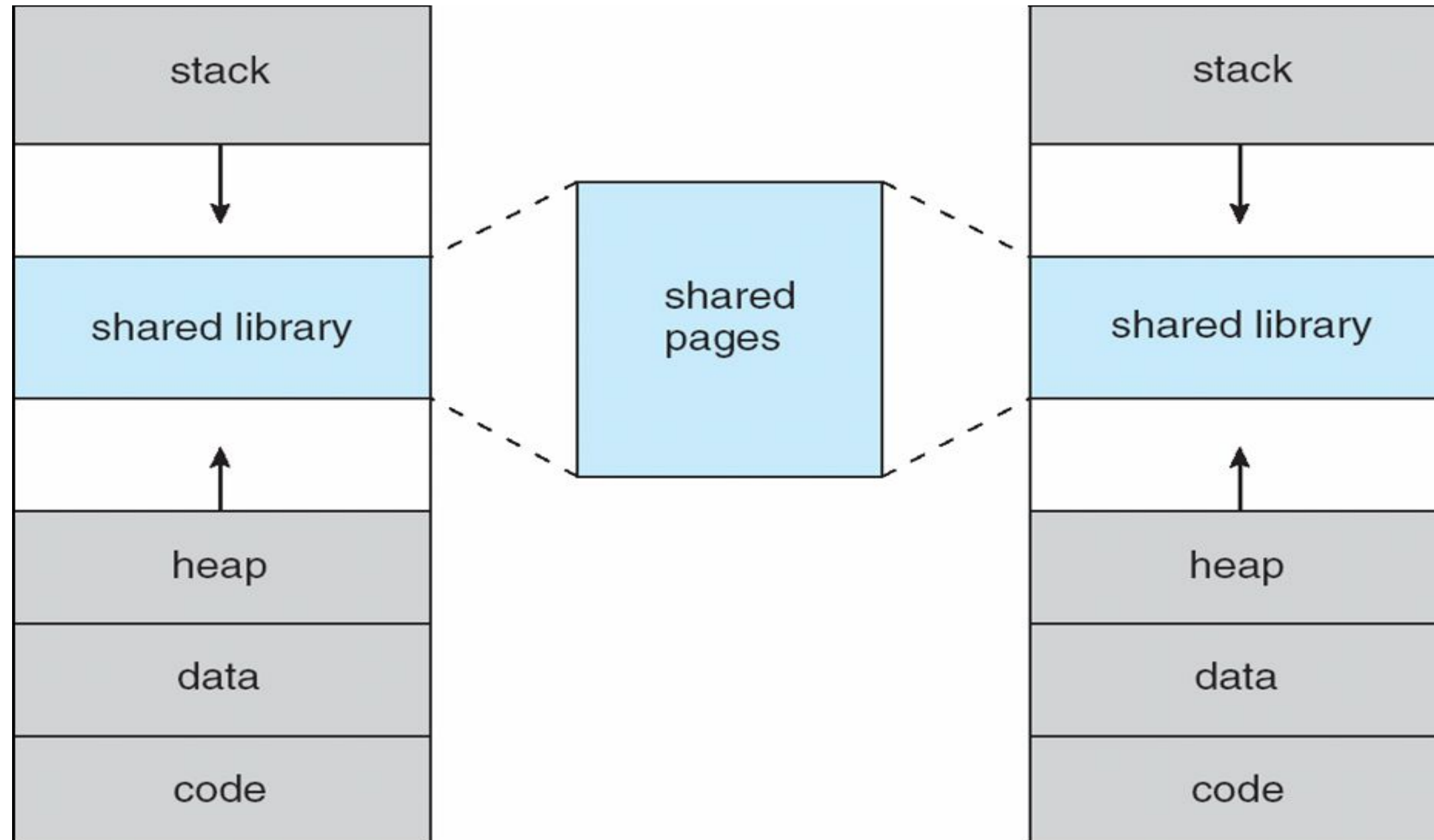
---

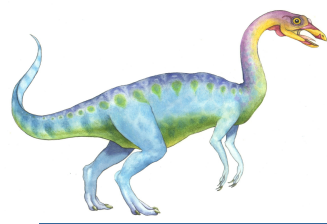
- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc
- System libraries shared via mapping into virtual address space





# Shared Library Using Virtual Memory





# Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed (Demand Paging)
  - Less I/O needed, no unnecessary I/O
  - 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**







# Valid-Invalid Entry

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

valid-invalid entry

Frame #	
	<b>v</b>
	<b>v</b>
	<b>v</b>
	<b>i</b>
....	
	<b>i</b>
	<b>i</b>
	<b>i</b>

During address translation, if valid–invalid entry in page table 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

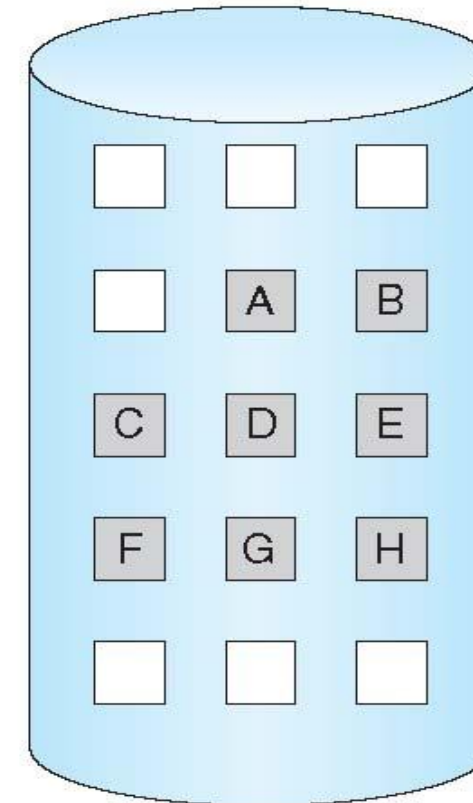
logical  
memory

valid-invalid bit		
frame		
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	A
5	
6	C
7	
8	
9	F
10	
11	
12	
13	
14	
15	

physical 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 via scheduled disk operation
4. Reset tables to indicate page now in memory  
Set validation bit = **v**
5. Restart the instruction that caused the page fault





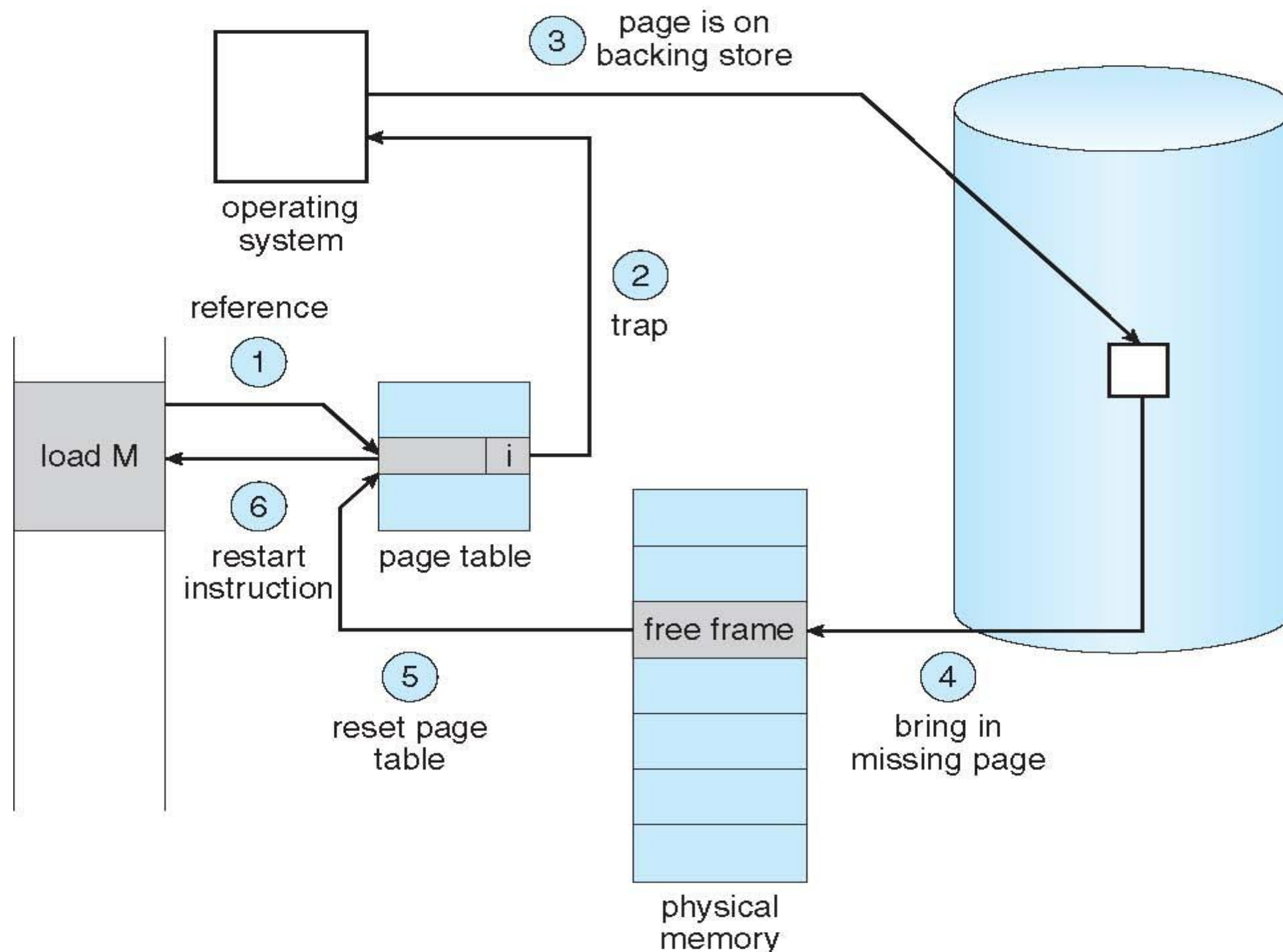
# Aspects of Demand Paging

- 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
  - Pain decreased because of **locality of reference**
- Hardware support needed for demand paging
  - Page table with valid / invalid bit
  - Secondary memory (swap device with **swap space**)
  - Instruction restarted





# Steps in Handling a Page Fault

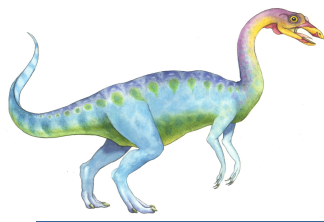




# Performance of Demand Paging

- Stages in Demand Paging
  1. Trap to the operating system
  2. Save the user registers and process state
  3. Determine that the interrupt was a page fault
  4. Check that the page reference was legal and determine the location of the page on the disk
  5. Find a free frame
  6. Issue a read from the disk to the free frame:
    1. Wait in a queue for this device until the read request is serviced
    2. Wait for the device seek and/or latency time
    3. Begin the transfer of the page to the free frame





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 (Cont.)

- Page Fault Rate  $0 \leq p \leq 1$ 
  - if  $p = 0$  no page faults
  - if  $p = 1$ , every reference is a fault

- Effective Access Time (EAT)

EAT =  $(1 - p)$  x memory access

+  $p$  (page fault overhead

+ swap page out

+ swap page in

+ restart overhead

)





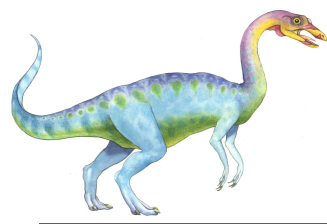


# Demand Paging Example

- Memory access time = 200 nanoseconds (1 nanosecond =  $10^{-9}$  second)
- Average page-fault service time = 8 milliseconds (=  $8 * 10^{-3}$  second)
- $$\begin{aligned} \text{EAT} &= (1 - p) \times 200 + p (8 \text{ milliseconds}) \\ &= (1 - p) \times 200 + p \times 8,000,000 \\ &= 200 + p \times 7,999,800 \end{aligned}$$
- If one access out of 1,000 causes a page fault, then  
$$\text{EAT} = 8.2 \text{ microseconds } (= 8.2 * 10^{-6} \text{ second}).$$

This is a slowdown by a factor of 40!!
- If want performance degradation < 10 percent. So EAT should be less than  $200 + 10\%$  of 200 (=EAT < 220)
  - $220 > 200 + 7,999,800 \times p$   
 $20 > 7,999,800 \times p$
  - $p < .0000025$
  - < one page fault in every 400,000 memory accesses



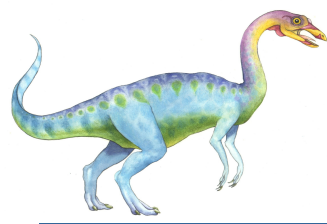


# Demand Paging Optimizations

---

- 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
- Data Pages allocated space on swap space when first paged out.
- On page fault pre-fetch other “nearby” pages





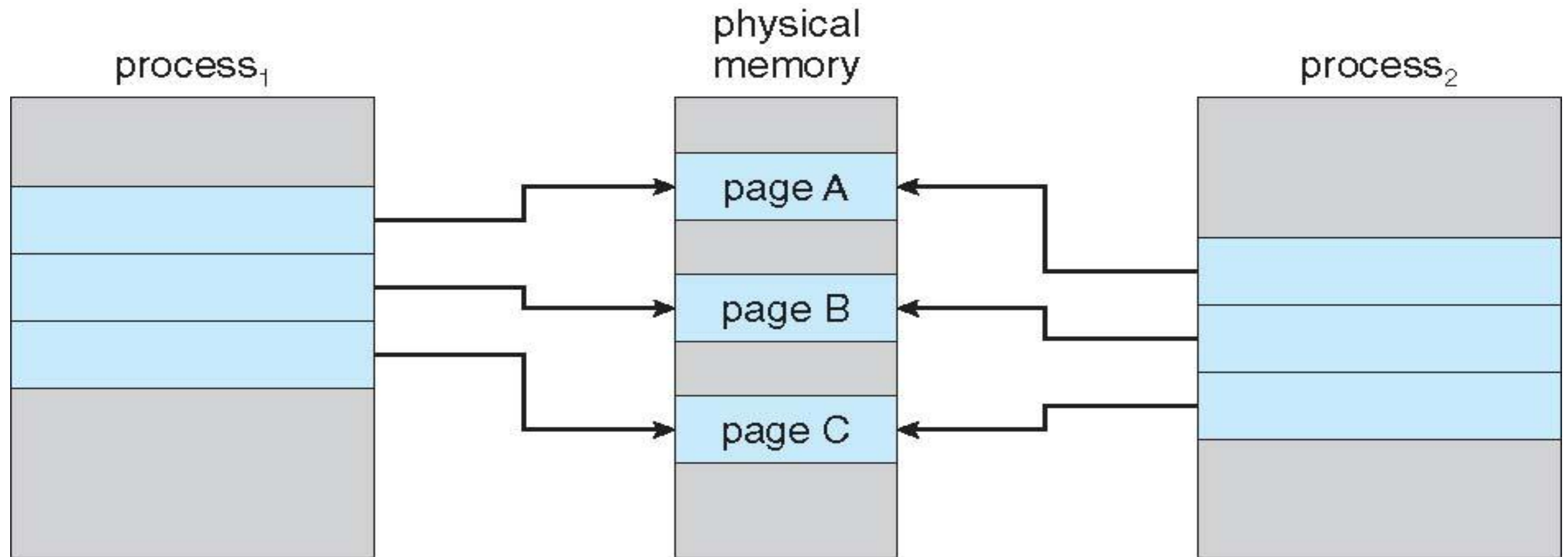
# Copy-on-Write

- **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
- `vfork()` variation on `fork()` system call has parent suspend and child using copy-on-write address space of parent
  - Designed to have child call `exec()`
  - Very efficient



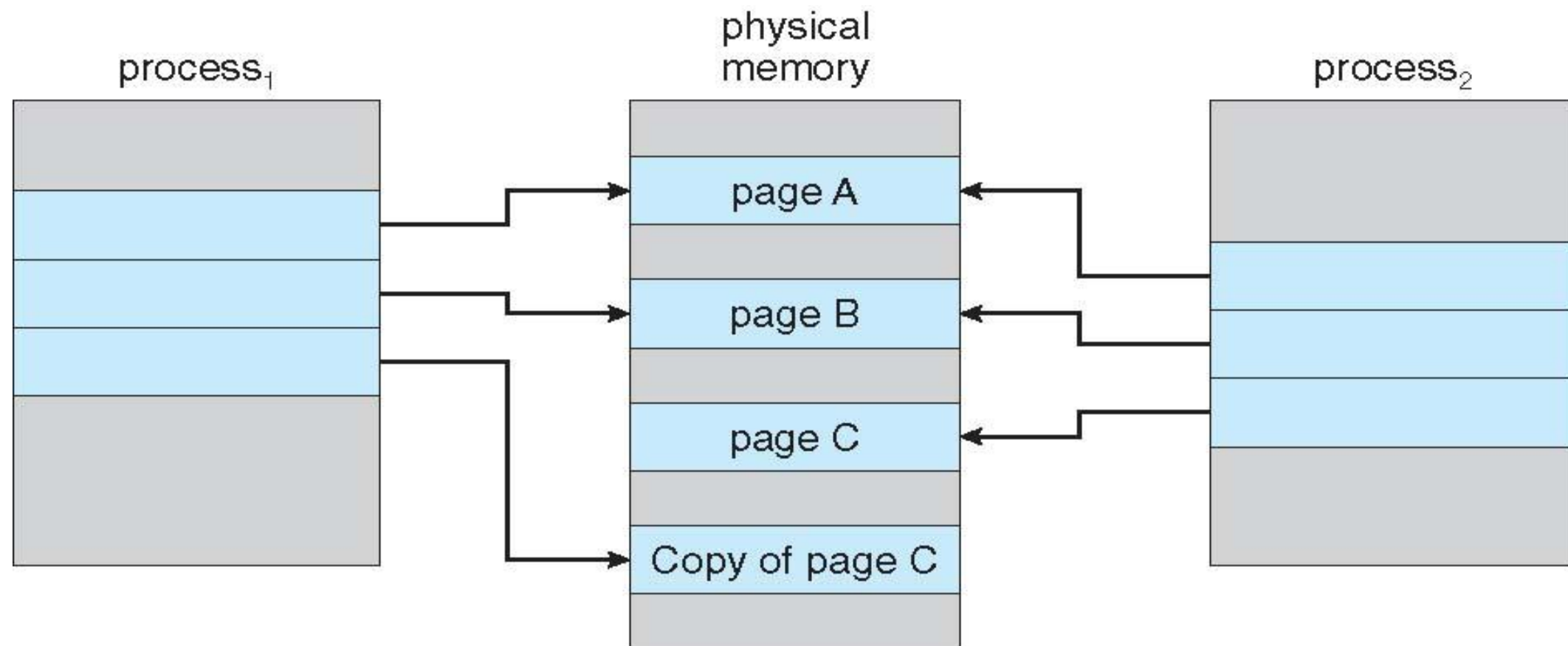


# Before Process 1 Modifies Page C





# After Process 1 Modifies Page C





# What Happens if There is no Free Frame?

- Used up by process pages
- Also in demand from the kernel, I/O buffers, etc
- How much to allocate to each?
- Page replacement – find some page in memory, but not really in use, page it out
  - Algorithm – terminate? swap out? replace the page?
  - Performance – want an algorithm which will result in minimum number of page faults
- Same page may be brought into memory several times





# Page Replacement

---

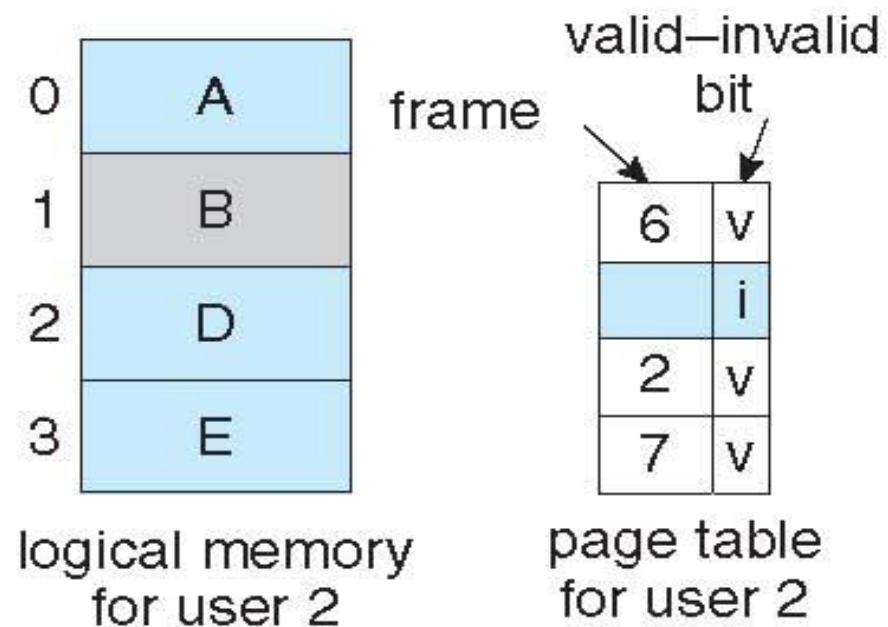
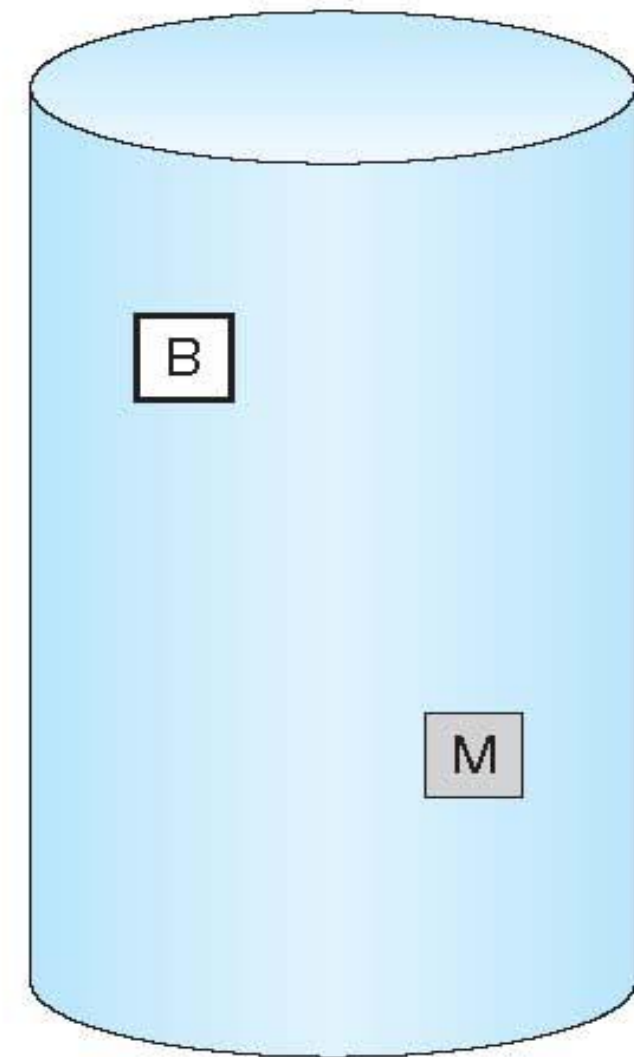
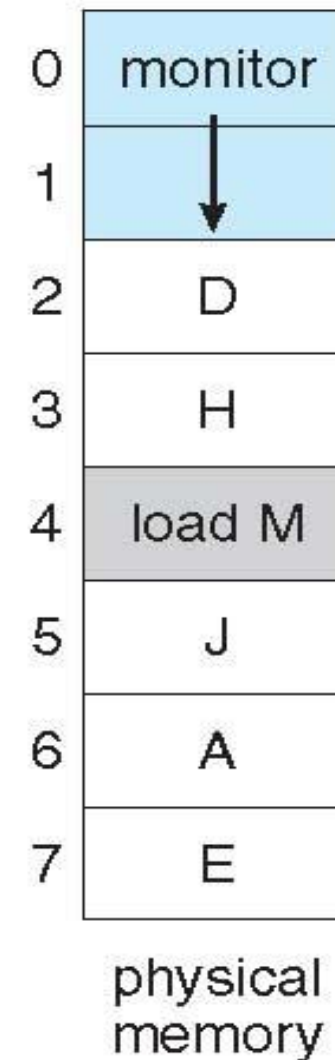
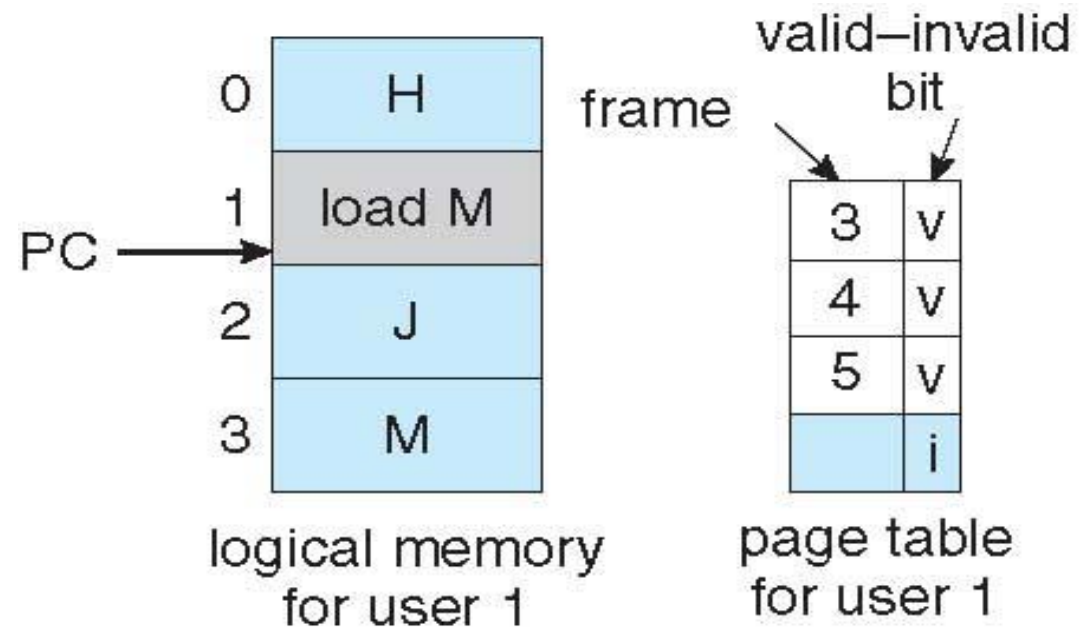
- Prevent over-allocation of memory by modifying page-fault service routine to include page replacement
- Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk







# Need For Page Replacement







# Basic Page Replacement

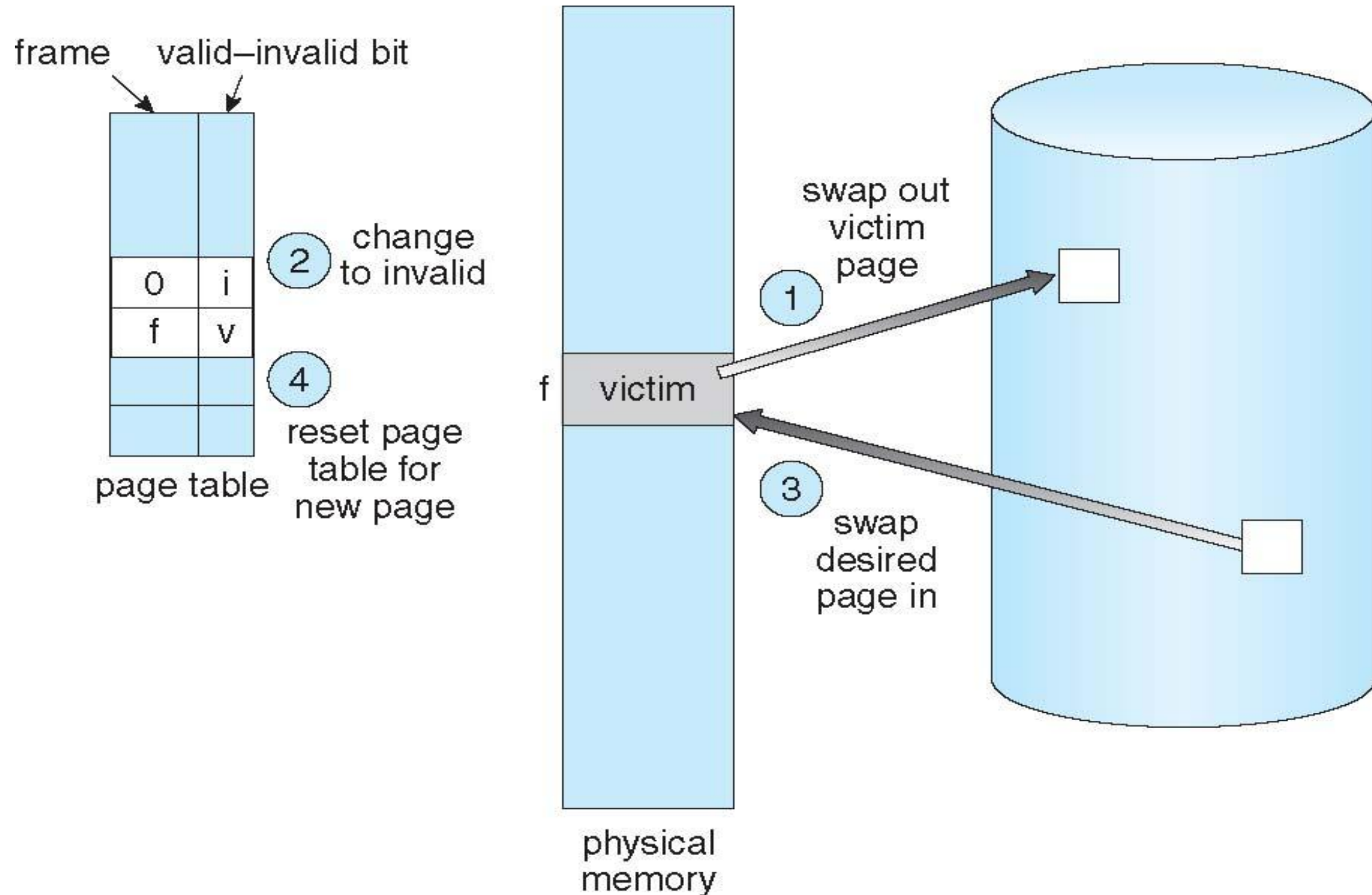
1. Find the location of the desired page on disk
2. Find a free frame:
  - If there is a free frame, use it
  - If there is no free frame, use a page replacement algorithm to select a **victim frame**
    - Write victim frame to disk if dirty
3. Bring the desired page into the (newly) free frame; update the page and frame tables
4. Continue the process by restarting the instruction that caused the trap

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





# Page Replacement





# Page and Frame Replacement Algorithms

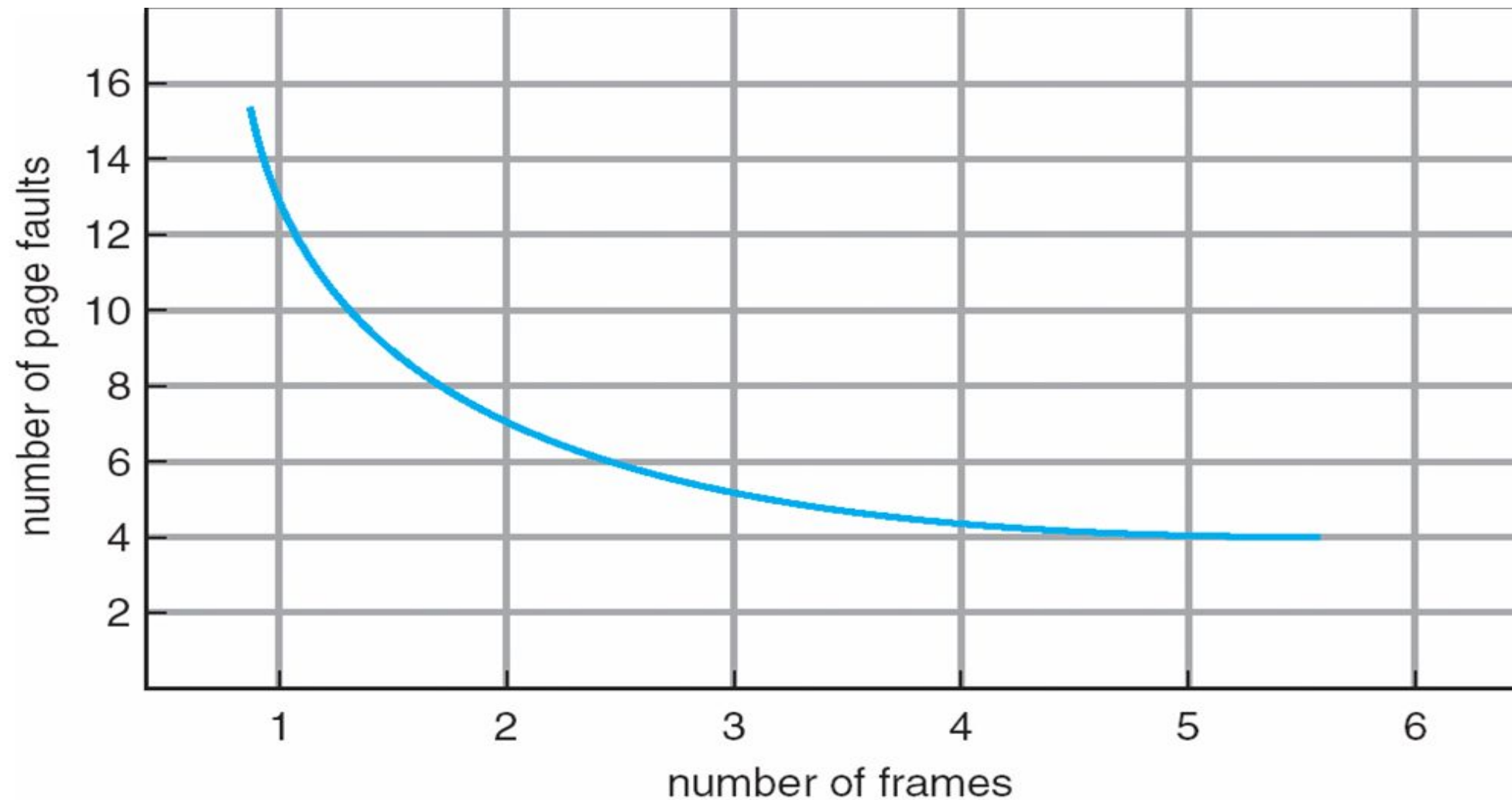
- **Frame-allocation algorithm** determines
  - How many frames to give each process
  - Which frames to replace
- **Page-replacement algorithm**
  - Want lowest page-fault rate on both first access and re-access
- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
  - String is just page numbers, not full addresses
  - Repeated access to the same page does not cause a page fault
- In all our examples, the reference string is

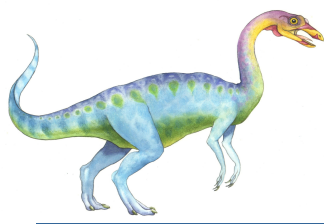
**7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**





# Graph of Page Faults Versus The Number of Frames





# First-In-First-Out (FIFO) Algorithm

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

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

15 page faults

- Can 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!

## 4 Belady's Anomaly

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





# FIFO Page Replacement

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																

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

0	0																		
1	1																		
3	2																		

7	7	7																	
1	0	0																	
2	2	1																	

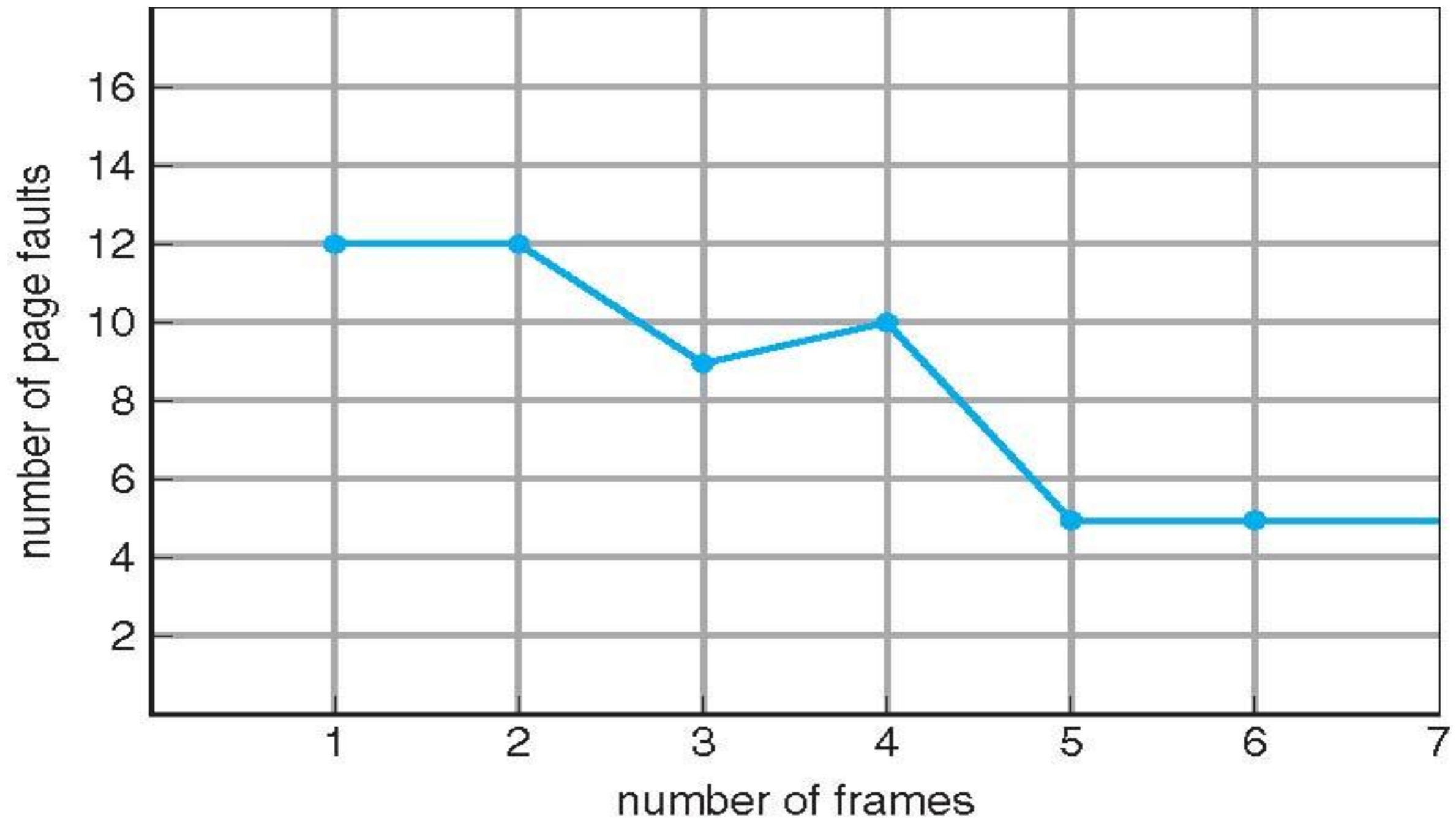
page frames

15 Faults





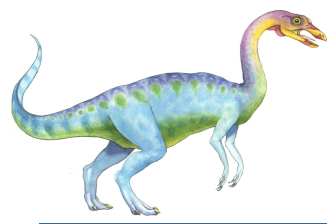
# FIFO Illustrating Belady's Anomaly



**Reference String is 123412512345**







# Optimal Algorithm

---

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



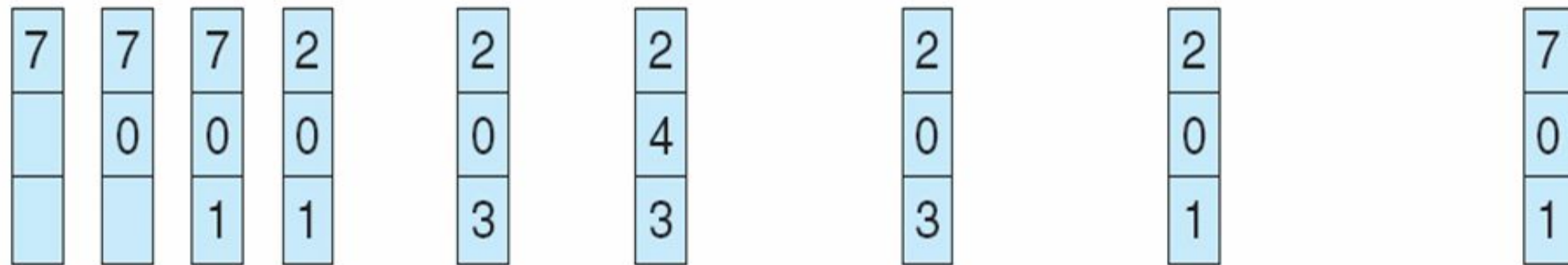




# Optimal Page Replacement

reference string

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



page frames

9 Faults





# 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

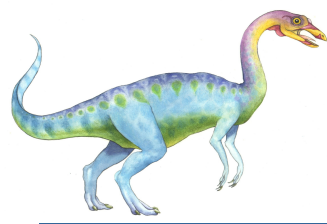
7	7	7	2		2		4	4	4	0		1		1		1
	0	0	0		0		0	0	3	3		3		0		0
		1	1		3		3	2	2	2		2		2		7

page frames

- 12 faults

LRU and OPT are cases of **stack algorithms** that don't have Belady's Anomaly





# 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
- Stack implementation
  - Keep a stack of page numbers in a double linked form:
  - Page referenced:
    - 4 move it to the top
    - 4 requires 6 pointers to be changed in worst case(for 5 frames)
- LRU and OPT are cases of **stack algorithms** that don't have Belady's Anomaly

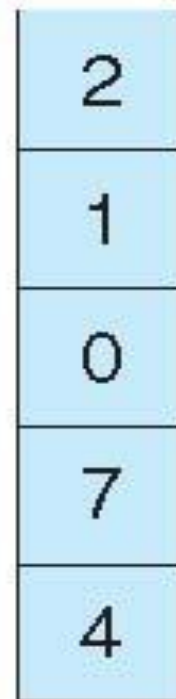




# Use Of A Stack To Record The Most Recent Page References

reference string

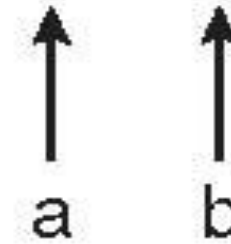
4 7 0 7 1 0 1 2 1 2 7 1 2



stack  
before  
a



stack  
after  
b

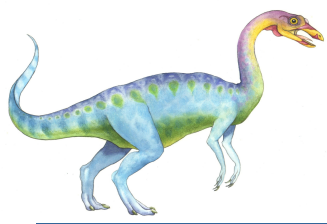




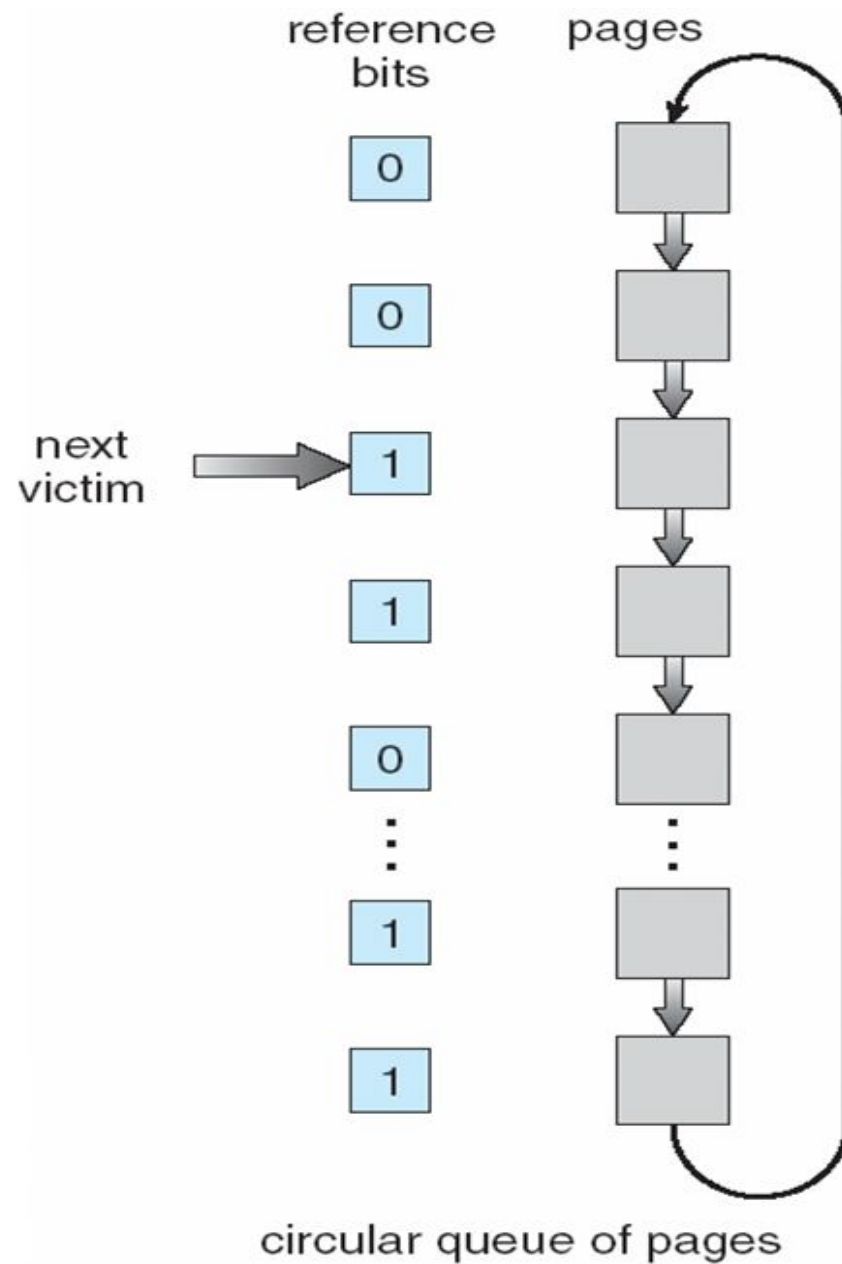
# LRU Approximation Algorithms

- LRU needs special hardware and still slow
- **Reference bit**
  - With each page associate a bit, initially = 0
  - When page is referenced bit set to 1
  - Replace any with reference bit = 0 (if one exists)
    - 4 We do not know the order, however
- **Second-chance algorithm**
  - Generally FIFO, plus hardware-provided reference bit
  - Clock replacement
  - If page to be replaced has
    - 4 Reference bit = 0  $\rightarrow$  replace it
    - 4 reference bit = 1 then:
      - set reference bit 0, leave page in memory
      - replace next page, subject to same rules

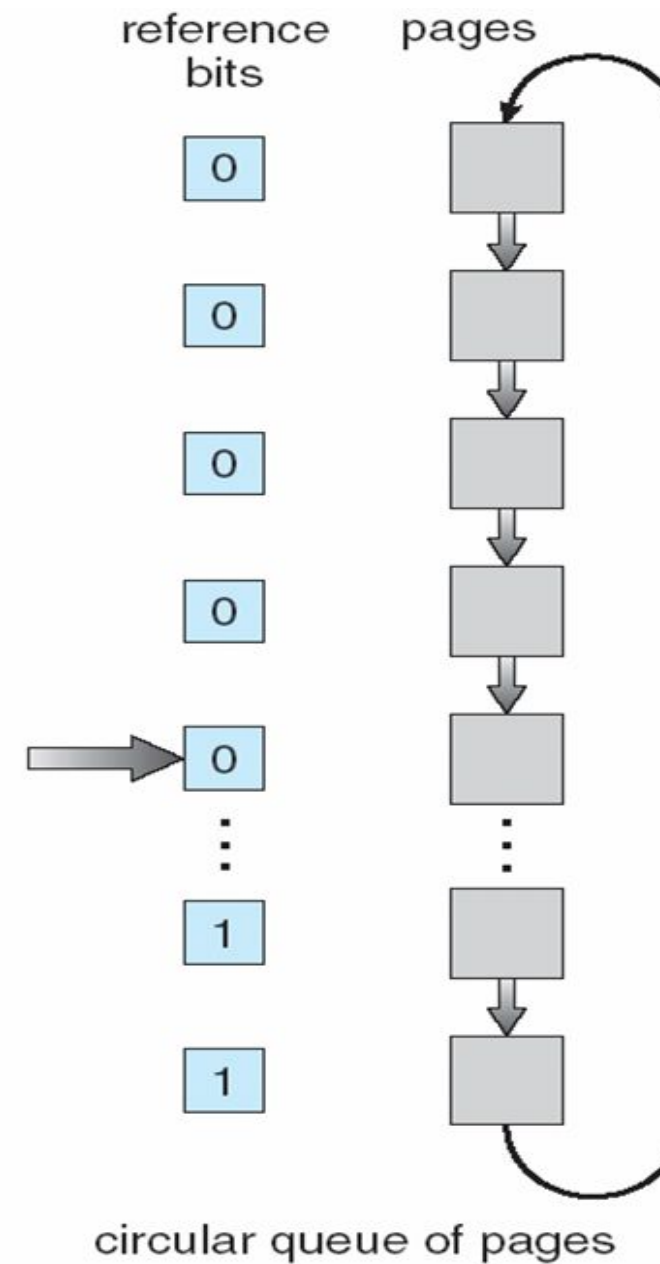




# Second-Chance (clock) Page-Replacement Algorithm



(a)



(b)





# Enhanced Second Chance with Dirty Bit

- The Second Chance Algorithm is also called a clock algorithm, with the pointer being the single clock hand.
- Use the dirty bit also in the decision making process. So we have a pair of bits (ref, dirty).
- Algorithm:
  - If victim is (0,0), replace.
  - If victim is (0,1), change to (0,0) and go to next page (remember dirty status elsewhere)
  - If victim is (1,0), change to (0,0) and go to next page.
  - If victim is (1,1), change to (0,1) and go to next page
- (0,1) and (1,0) get a second chance; (1,1) gets two chances







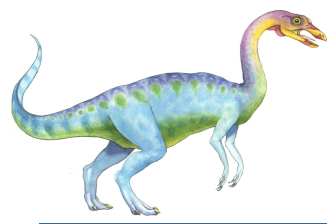
# Counting Algorithms

---

- Keep a counter of the number of references that have been made to each page
  - Not common
- **LFU (Least Frequently Used) Algorithm:** replaces page with smallest count
- **MFU (Most Frequently Used) Algorithm:** based on the argument that the page with the smallest count was probably just brought in and has yet to be used







# Allocation of Frames

---

- Each process needs a *minimum* number of frames
- *Maximum* of course is the total frames in the system or the size of the process
- Two major allocation schemes
  - fixed allocation
  - priority allocation
- Many variations
- Allocation and replacement algorithms are related





# Fixed Local 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
- Priority Allocation – allocate acc. to priority of processes





# 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





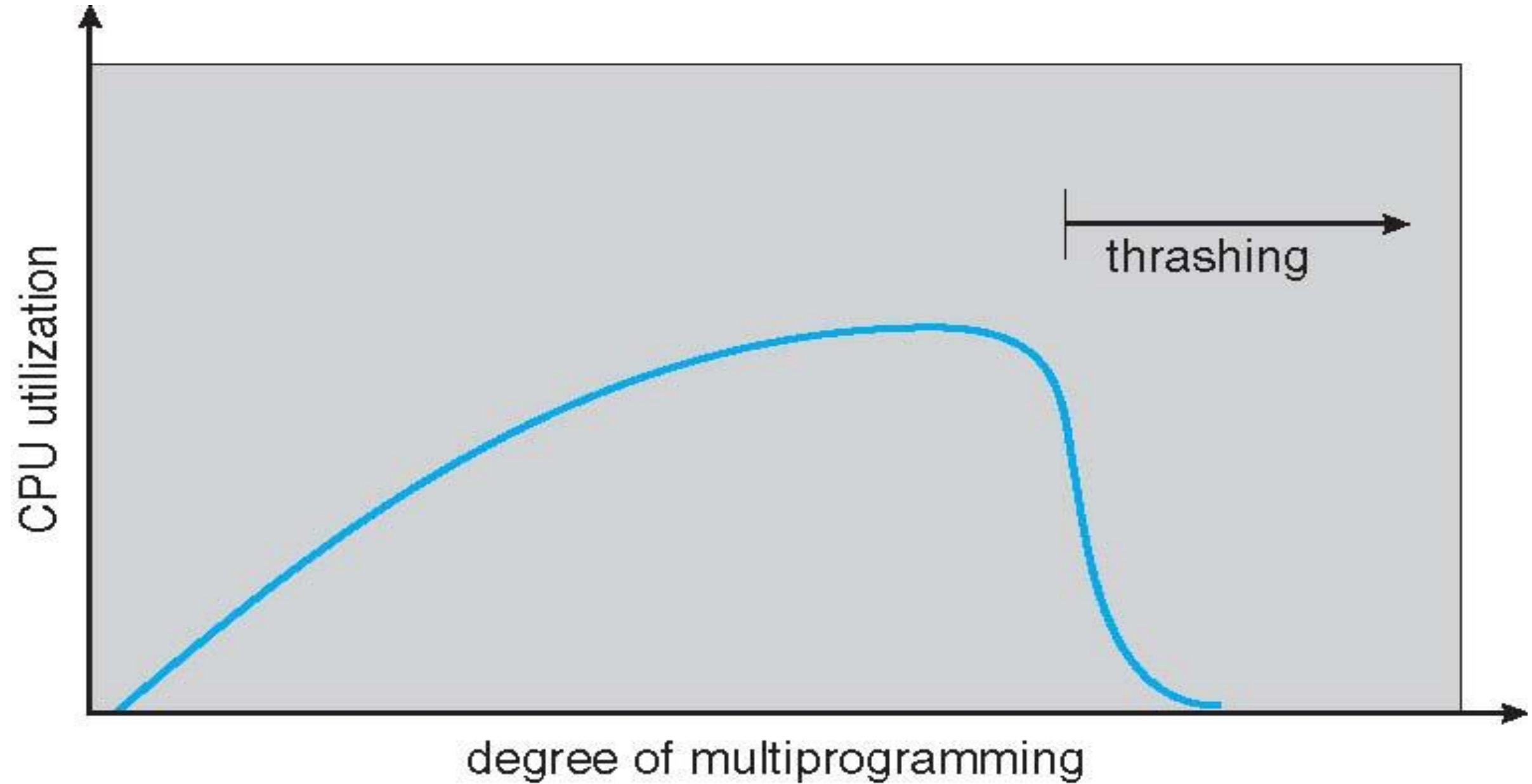
# Thrashing

- 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:
    - 4 Low CPU utilization
    - 4 Operating system thinking that it needs to increase the degree of multiprogramming
    - 4 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?

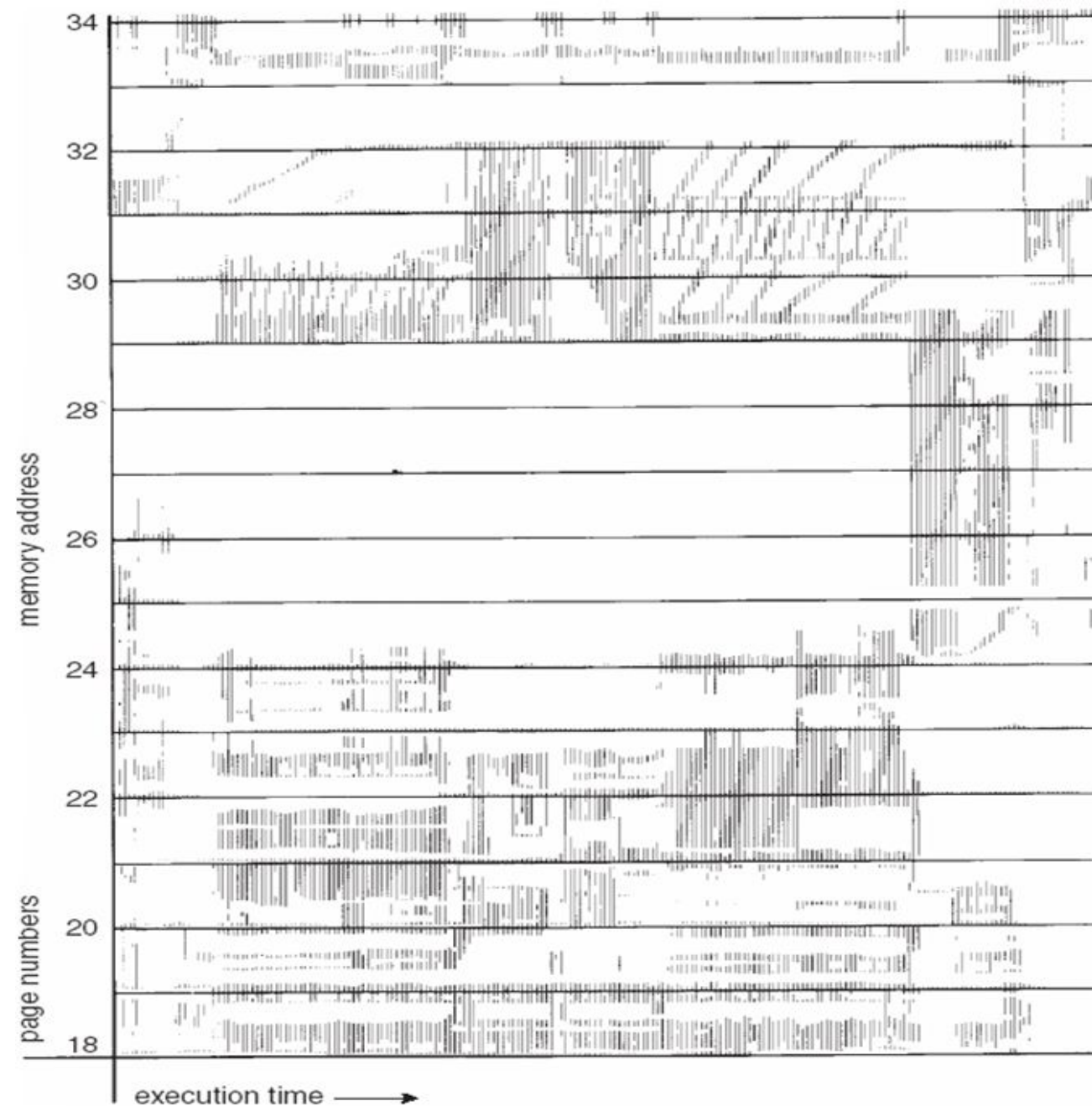
## Locality model

- Process migrates from one locality to another
  - Localities may overlap
- 
- Why does thrashing occur?  
 $\Sigma$  size of locality  $>$  size of memory allocated
    - 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  
Example: 10,000 instructions
- $WSS_i$  (working set of Process  $P_i$ ) =  
total number of pages referenced in the most recent  $\Delta$  (varies in time)
  - if  $\Delta$  too small will not encompass entire locality
  - if  $\Delta$  too large will encompass several localities
  - if  $\Delta = \infty \Rightarrow$  will encompass entire program
- $D = \Sigma WSS_i \equiv$  total demand frames
  - Approximation of locality
- if  $D > m \Rightarrow$  Thrashing (**m is the number of frames allocated**)
- Policy if  $D > m$ , then suspend or swap out one of the processes



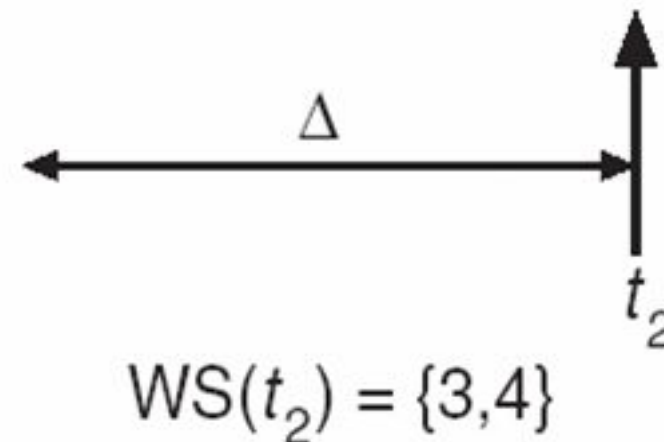
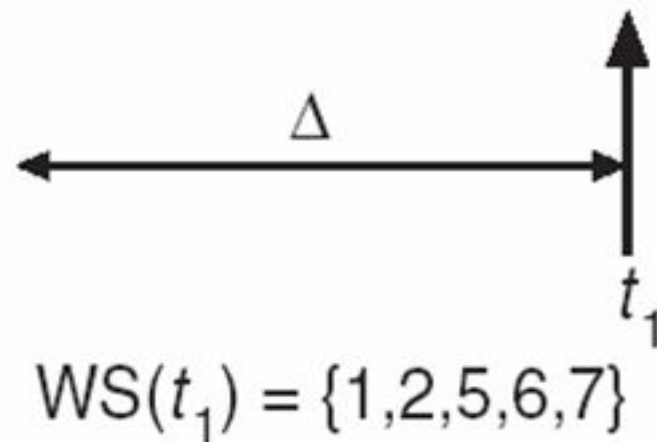


# Working-set model

- Working Set Window is 10 references:

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 ...



- After every reference, the Working Set may change
- Too expensive to implement





# Keeping Track of the Working Set

---

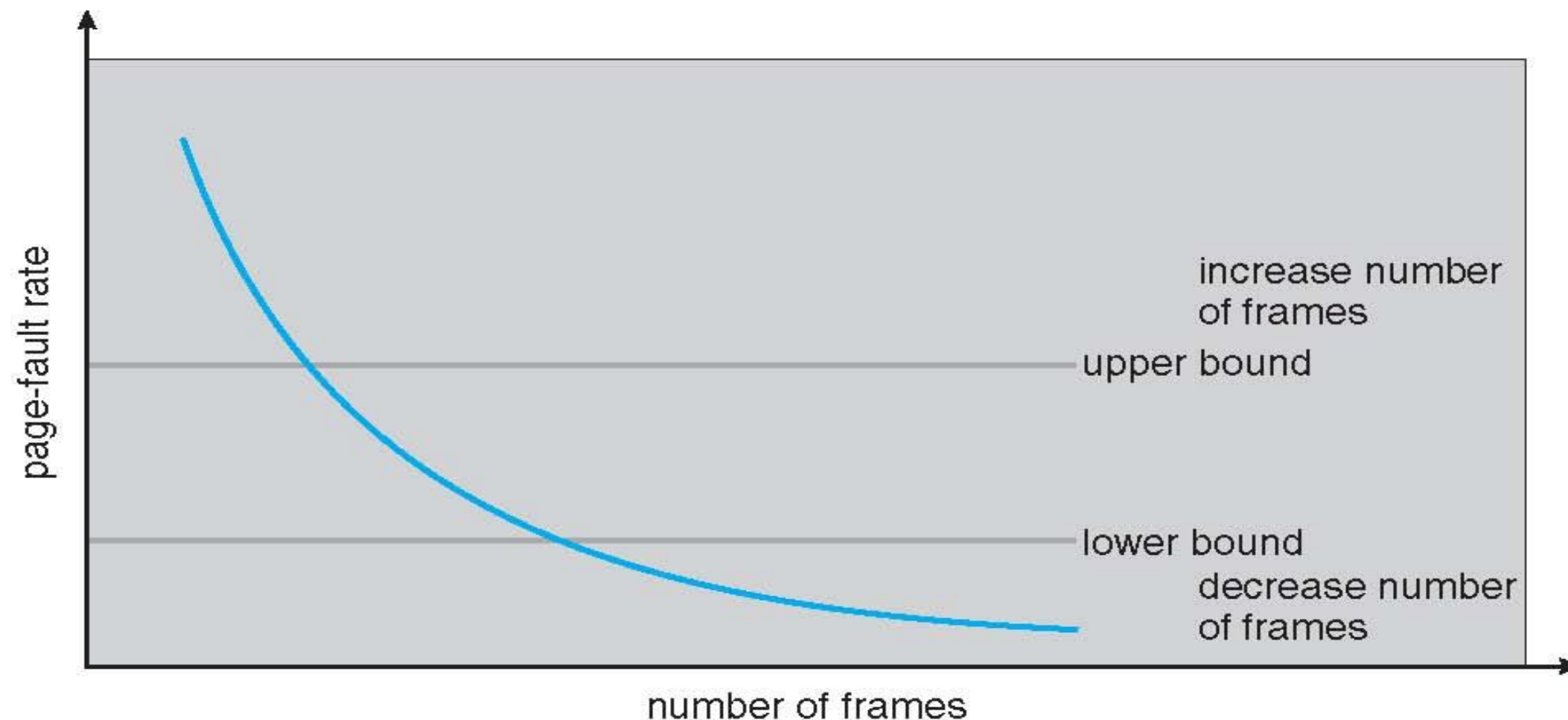
- Approximate with interval timer + a reference bit
- Example:  $\Delta = 10,000$ 
  - Timer interrupts after every 5000 time units
  - Keep in memory 2 bits for each page
  - Whenever a timer interrupts copy and sets the values of all reference bits to 0
  - If one of the bits in memory = 1  $\Rightarrow$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





# Page-Fault Frequency

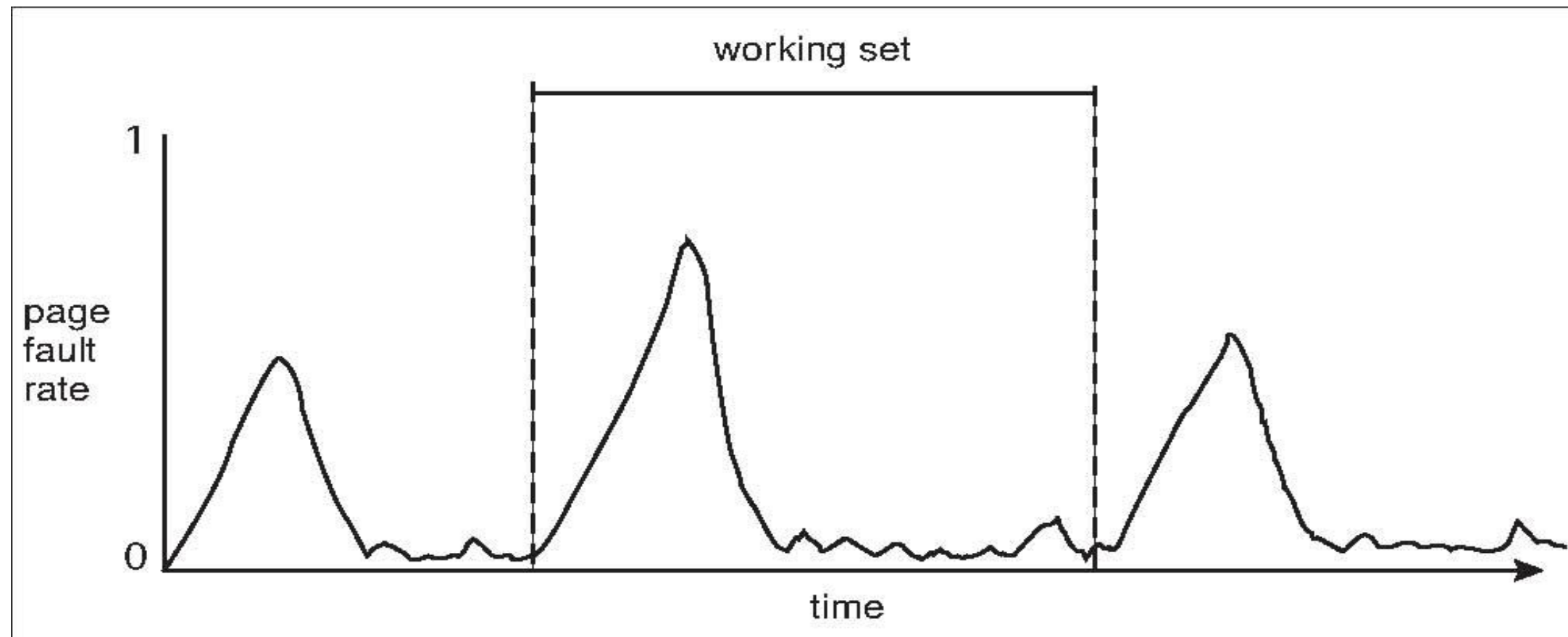
- More direct approach than WSS
- Establish “acceptable” **page-fault frequency** rate and use local replacement policy
  - If actual rate too low, process loses frame
  - If actual rate too high, process gains frame





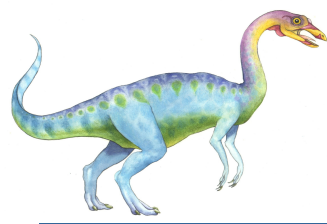
# Working Sets and Page Fault Rates

- Assume Working Sets fit in memory always.
- The page fault rates then are seen to have the following pattern:



- Peaks represent a move from one locality to another.





# 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 and speeds file access by driving file I/O through memory rather than `read()` and `write()` system calls
- Also allows several processes to map the same file allowing the pages in memory to be shared
- But when does written data make it to disk?
  - Periodically and / or at file `close()` time
  - For example, when the pager scans for dirty pages





# Memory-Mapped File Technique for all I/O

---

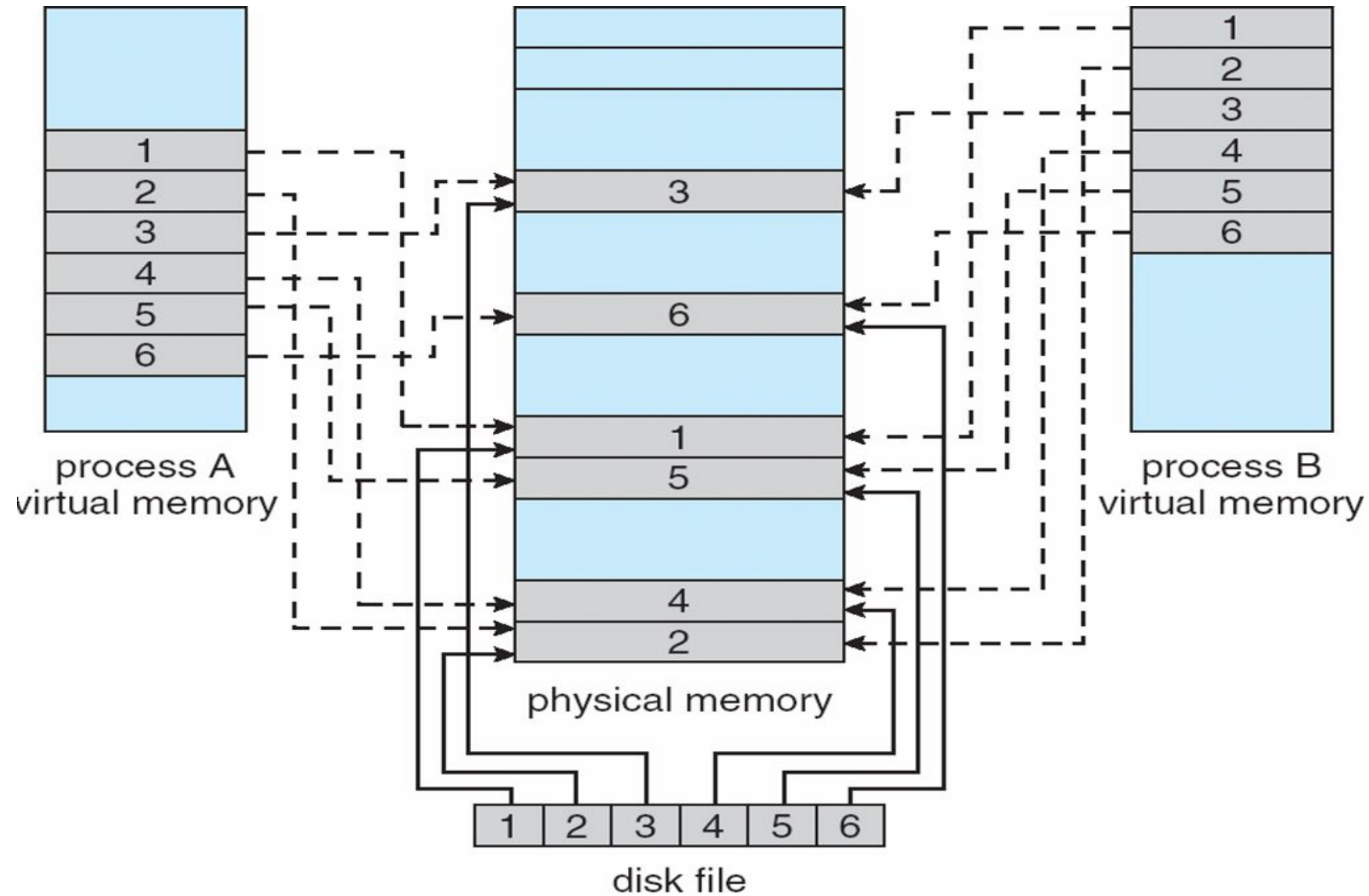
- Some OSes use memory mapped files for standard I/O
- Process can explicitly request memory mapping a file via `mmap ( )` system call
  - Now file mapped into process address space
- Memory mapped files can be used for shared memory (although again via separate system calls)







# Memory Mapped Files



# End of Chapter 9

---

