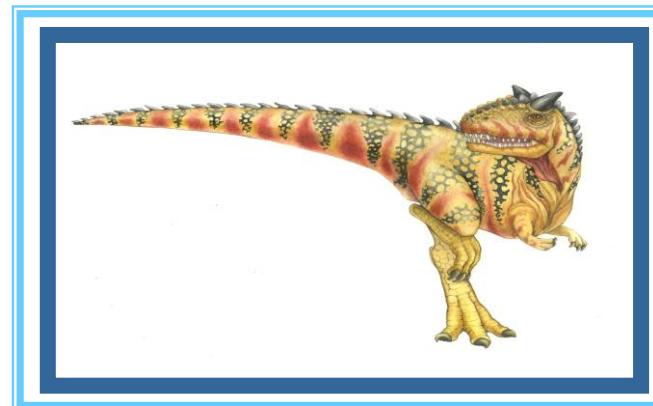
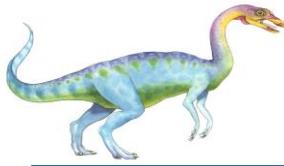


Chapter 9: Virtual Memory

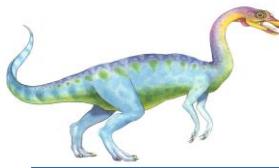




Chapter 9: Virtual Memory

- Background
- Demand Paging
- Copy-on-Write
- Page Replacement
- Allocation of Frames
- Thrashing
- Memory-Mapped Files
- Allocating Kernel Memory

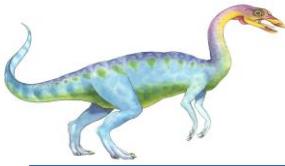




Timeline for Completion of OS Portion

No. of classes	Lecture No.	Date	Portion
3	L24, L25* & L26*	29-10-2022 Sat	Virtual Memory
1	L27 (CC*)	02-11-2022 Wed	Disk scheduling – 1
1	L28	03-11-2022 Thurs	Disk scheduling – 2
1	L29	05-11-2022 Sat	Disk scheduling – 3
1	L30	08-11-2022 Tues	Files
1	L31	10-11-2022 Thurs	Wrapping up

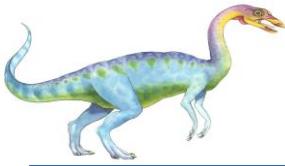




Objectives

- 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
- To examine the relationship between shared memory and memory-mapped files
- To explore how kernel memory is managed

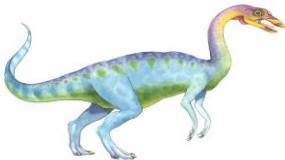




Background

- 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





Background (Cont.)

- **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
 - More programs running concurrently
 - Less I/O needed to load or swap processes





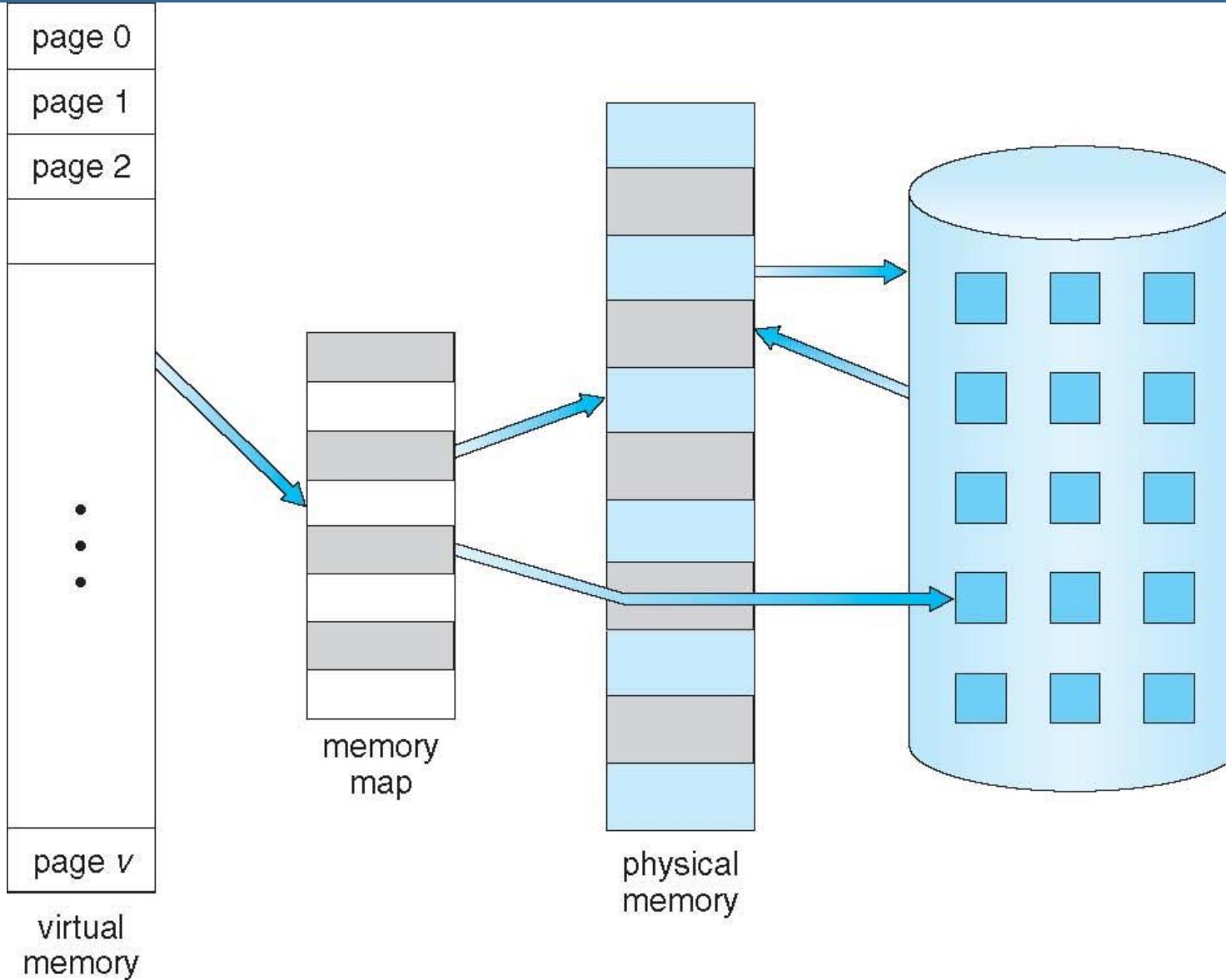
Background (Cont.)

- **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





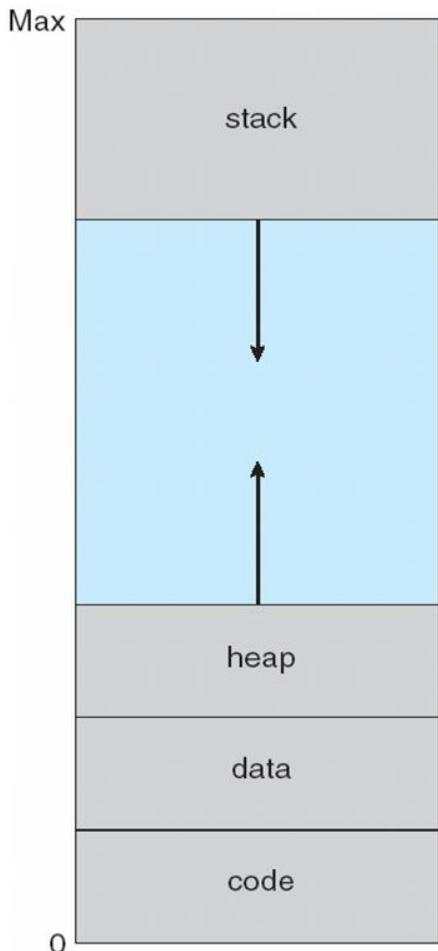
Virtual Memory That is Larger Than Physical Memory

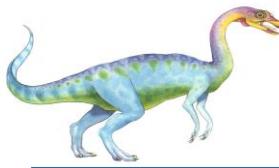




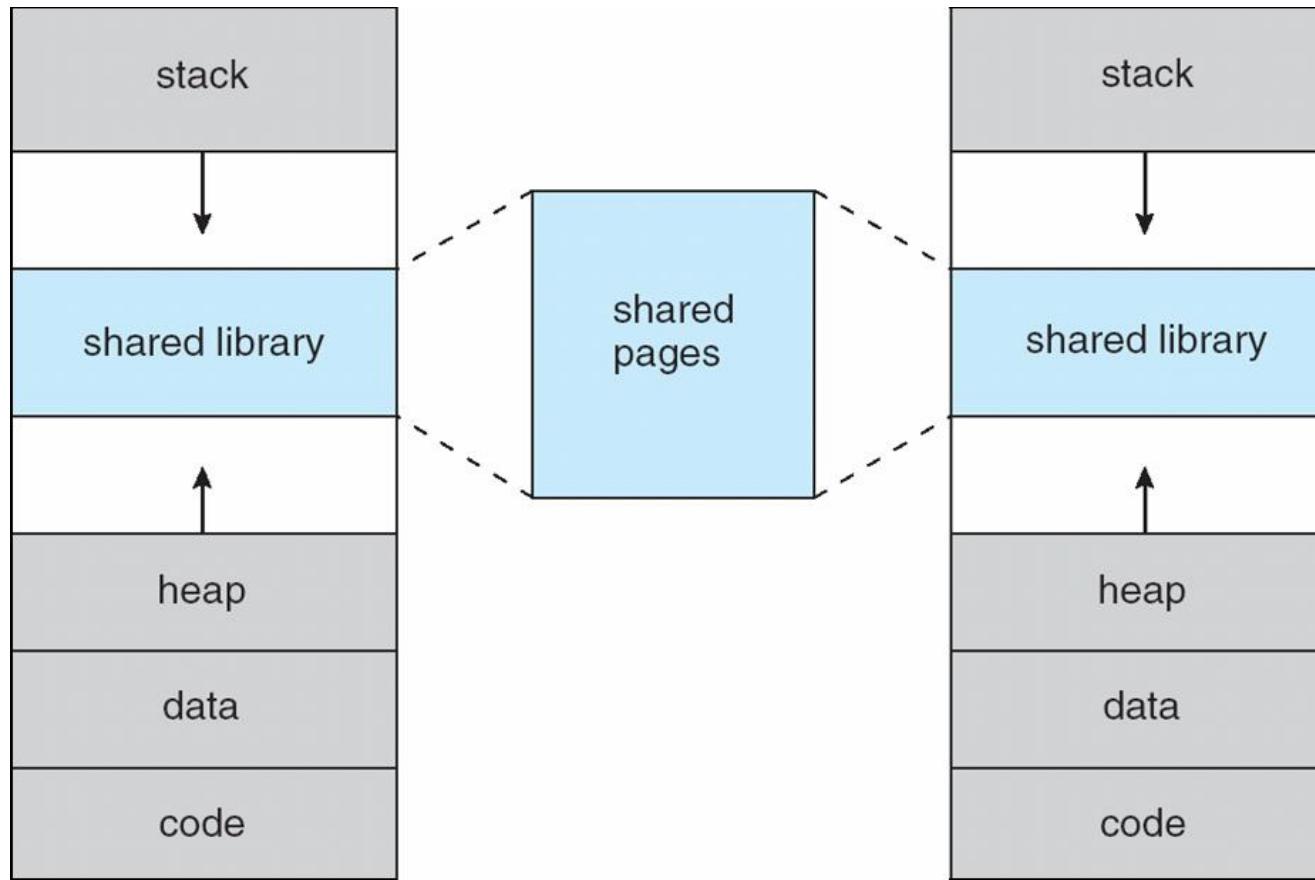
Virtual-address Space

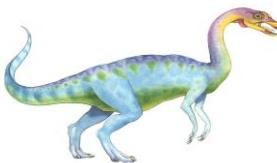
- Usually design logical address space for stack to start at Max logical address and grow “down” while heap grows “up”
 - Maximizes address space use
 - Unused address space between the two is hole
 - No physical memory needed until heap or stack grows to a given new page
- Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc
- System libraries shared via mapping into virtual address space
- Shared memory by mapping pages read-write into virtual address space
- Pages can be shared during `fork()`, speeding process creation





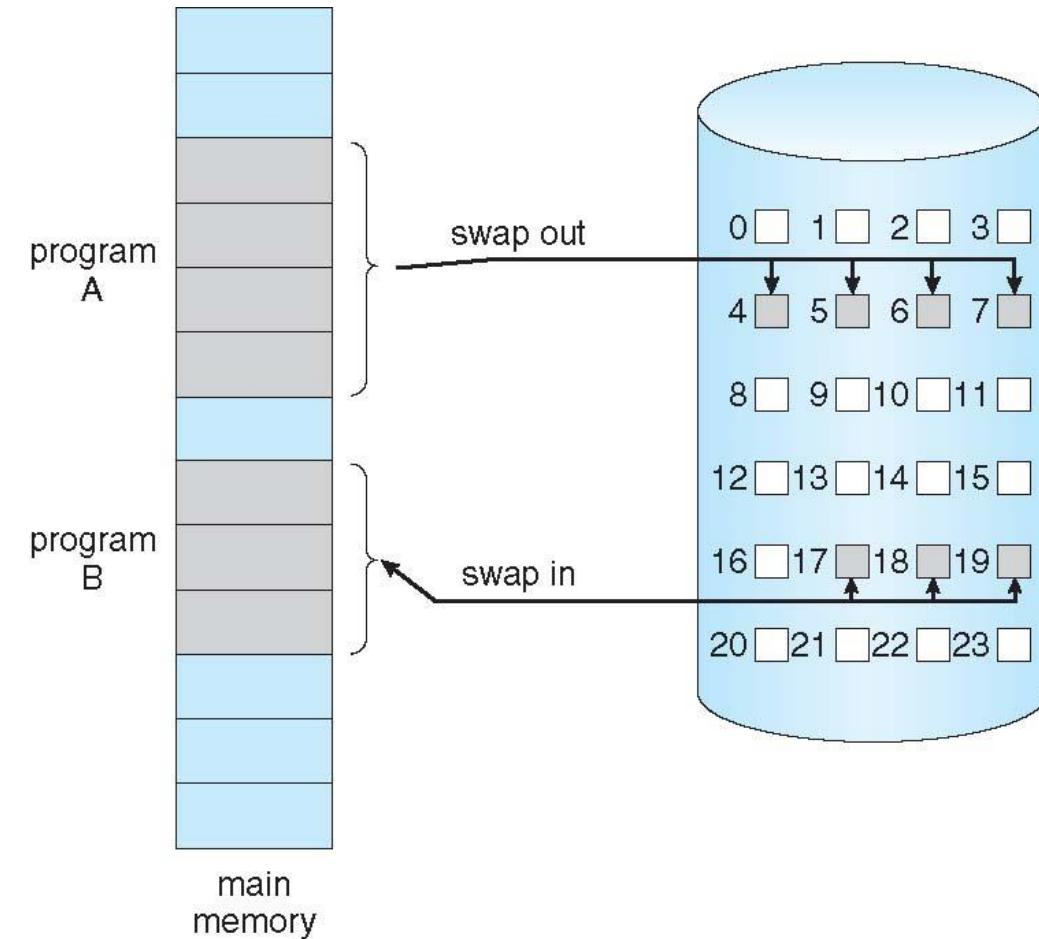
Shared Library Using Virtual Memory





Demand Paging

- Could bring entire process into memory at load time
- Or bring a page into memory only when it is needed
 - 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 ⇒ reference to it
 - invalid reference ⇒ abort
 - not-in-memory ⇒ bring to memory
- **Lazy swapper** – never swaps a page into memory unless page will be needed
 - Swapper that deals with pages is a **pager**





Basic Concepts

- With swapping, pager guesses which pages will be used before swapping out again
- Instead, pager brings in only those pages into memory
- How to determine that set of pages?
 - Need new MMU functionality to implement demand paging
- If pages needed are already **memory resident**
 - No difference from non demand-paging
- If page needed and not memory resident
 - Need to detect and load the page into memory from storage
 - ▶ Without changing program behavior
 - ▶ Without programmer needing to change code





Valid-Invalid Bit

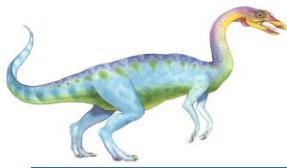
- With each page table entry a valid–invalid bit is associated (**v** ⇒ in-memory – **memory resident**, **i** ⇒ not-in-memory)
- Initially valid–invalid bit is set to **i** on all entries
- Example of a page table snapshot:

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** ⇒ page fault

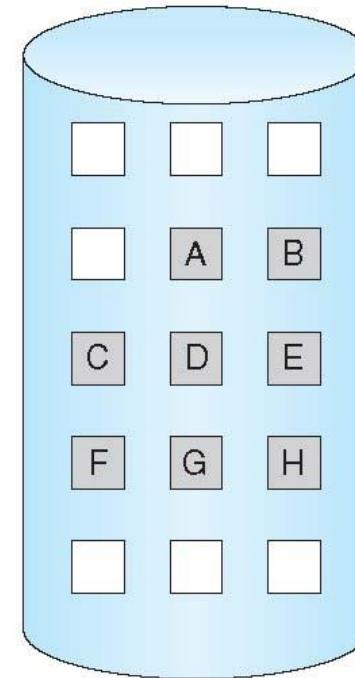
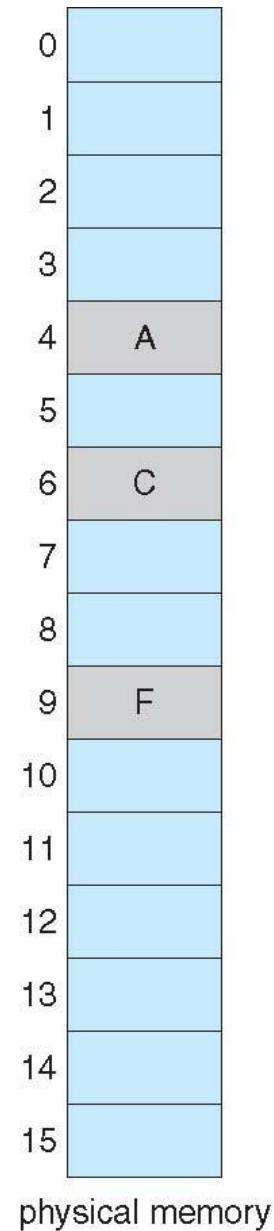
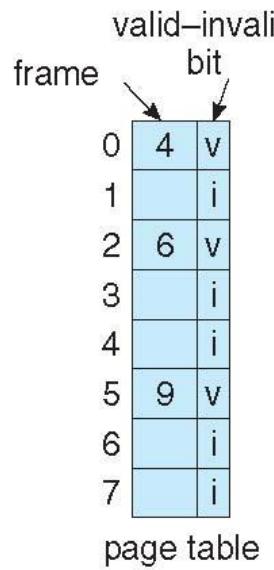




Page Table When Some Pages Are Not in Main Memory

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

logical memory





Page Fault

- If there is a reference to a page, first reference to that page will trap to operating system:

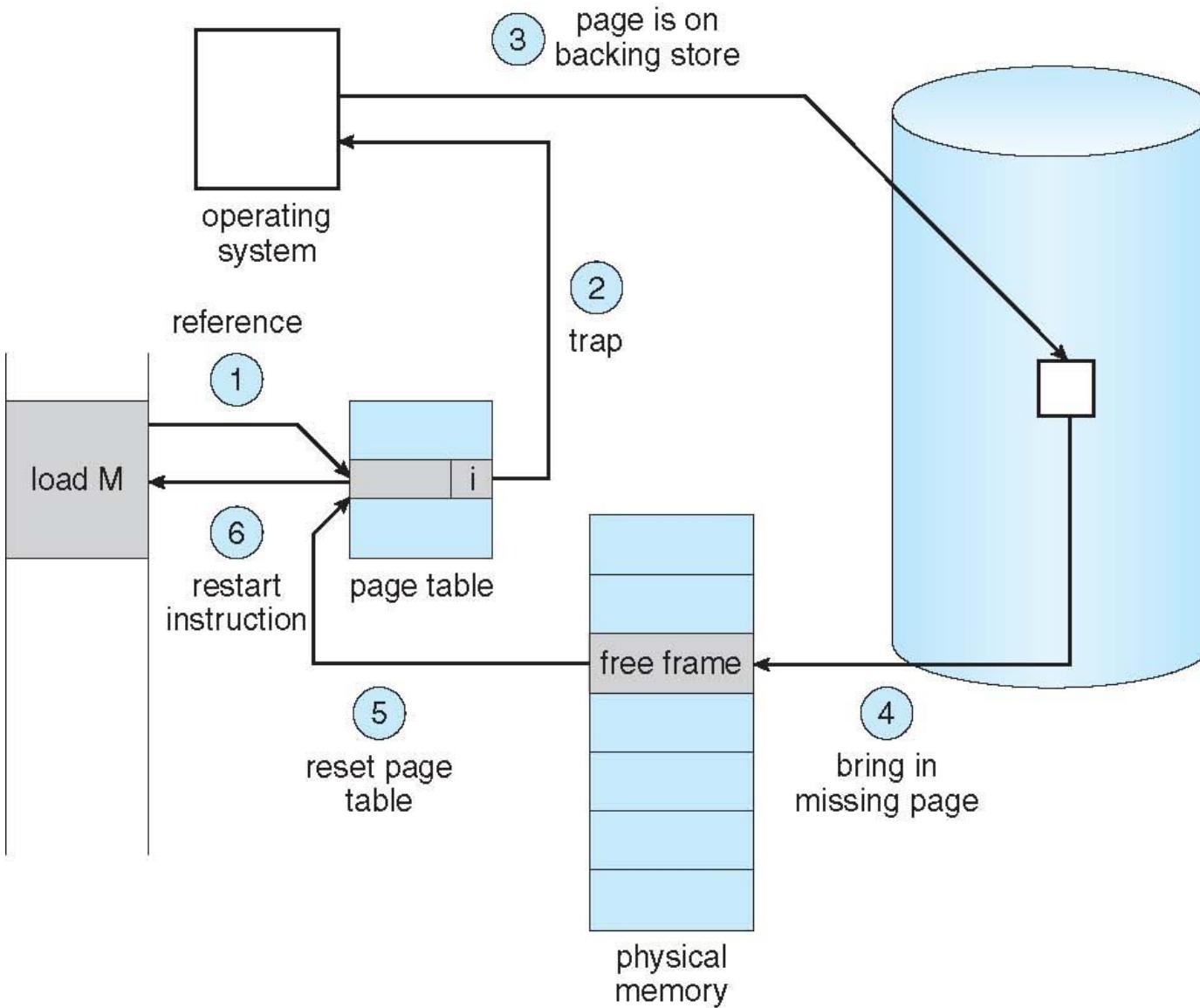
page fault

1. Operating system looks at another table to decide:
 - Invalid reference \Rightarrow abort
 - Just not in memory
2. Find free 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





Steps in Handling a Page Fault

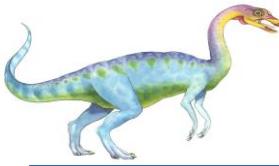




Aspects of Demand Paging

- Extreme case – start process with *no* pages in memory
 - OS sets instruction pointer to first instruction, 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**
- Hardware support needed for demand paging
 - Page table with valid / invalid bit
 - Secondary memory (swap device with **swap space**)
 - Instruction restart



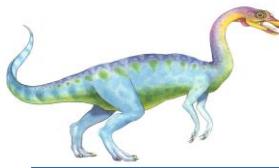


Performance of Demand Paging

□ Stages in Demand Paging (worse case)

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. Issue a read from the disk to a 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 a 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.)





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
 $EAT = 8.2 \text{ microseconds.}$
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





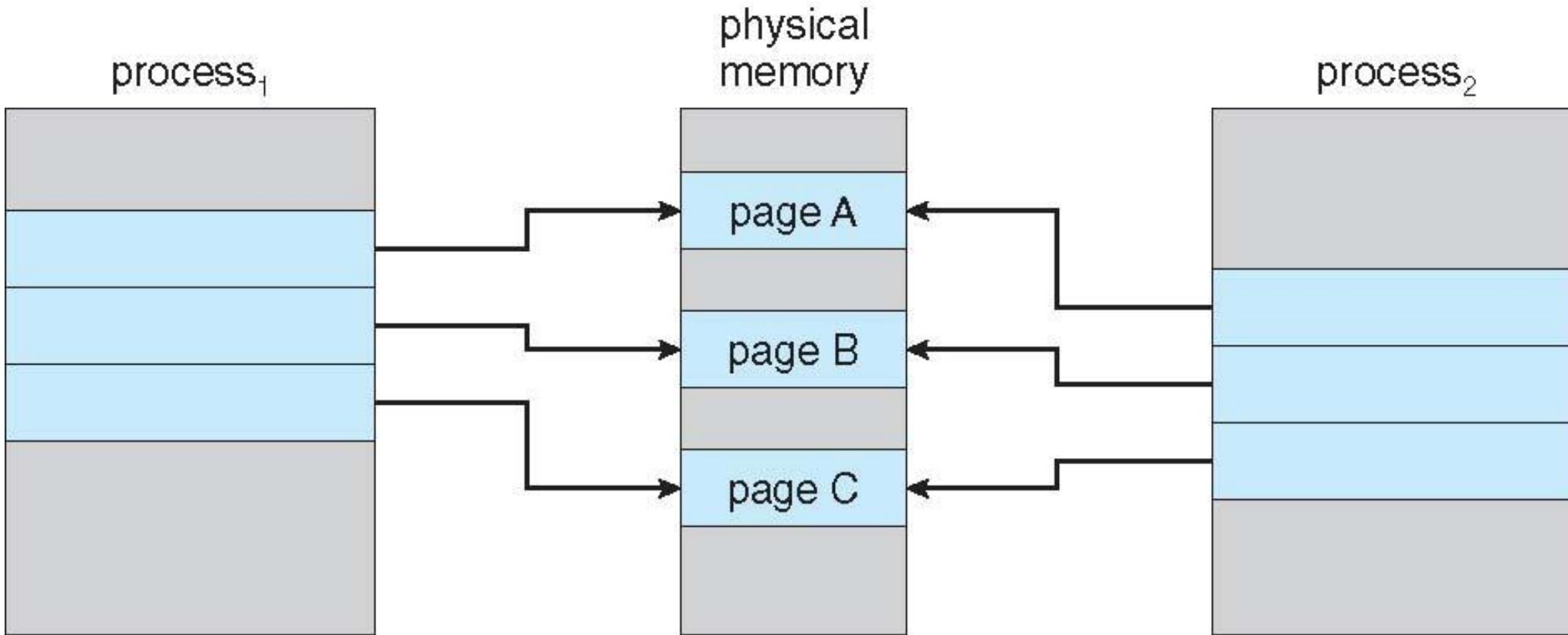
Copy-on-Write

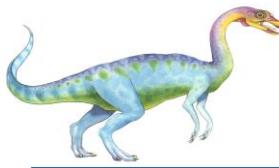
- **Copy-on-Write** (COW) allows both parent and child processes to initially **share** the same pages in memory
 - If either process modifies a shared page, only then is the page copied
- COW allows more efficient process creation as only modified pages are copied
- In general, free pages are allocated 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?
- `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



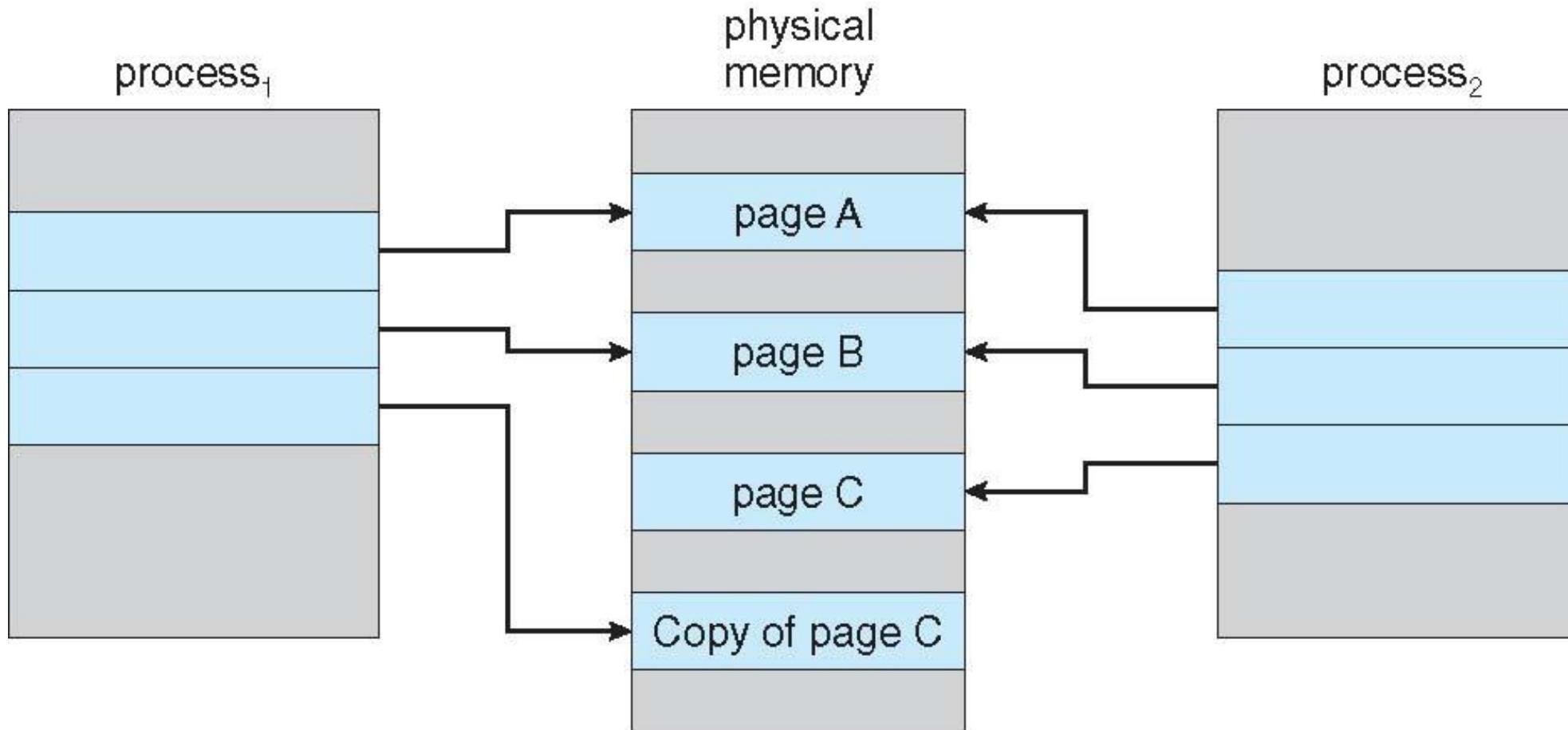


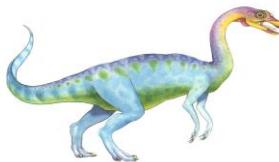
Before Process 1 Modifies Page C





After Process 1 Modifies Page C

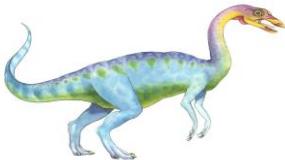




What Happens if There is no Free Frame?

- 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





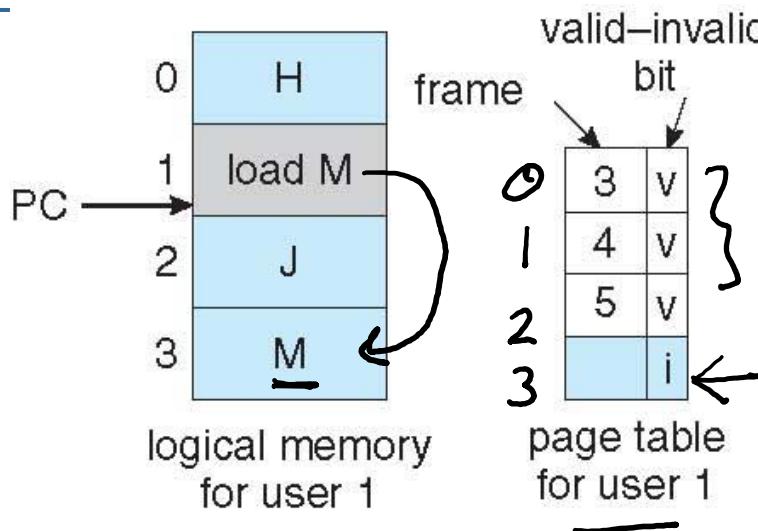
Page Replacement

- Prevent **over-allocation** of memory by modifying page-fault service routine to include page replacement
- Use **modify (dirty) bit** ~~=X~~ O 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



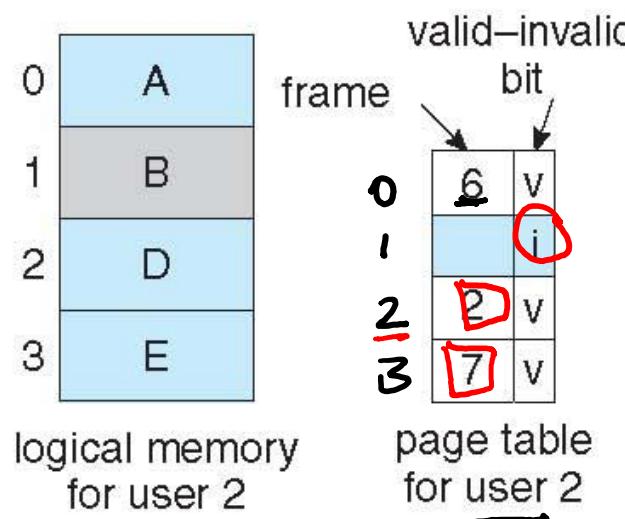


Need For Page Replacement



Opt 1²
Termination

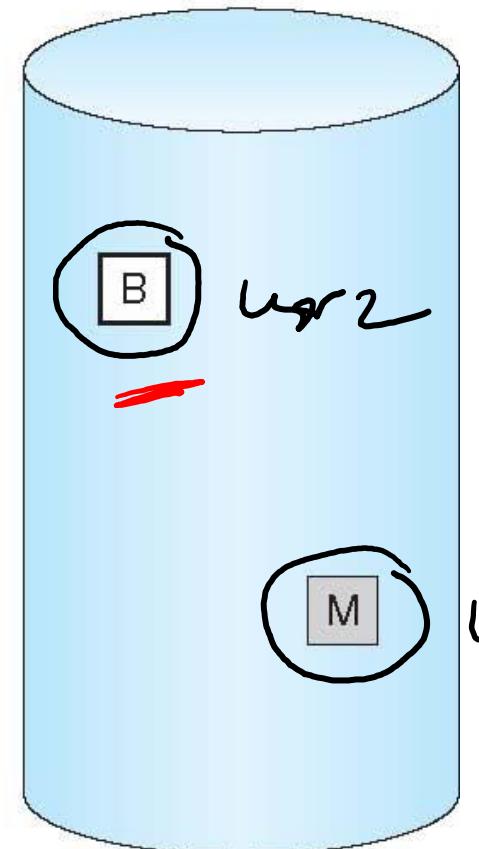
Opt 2.
Page Replacement



0	monitor
1	
2	D
3	H
4	load M
5	J
6	A
7	E

physical memory

frame used





Basic Page Replacement

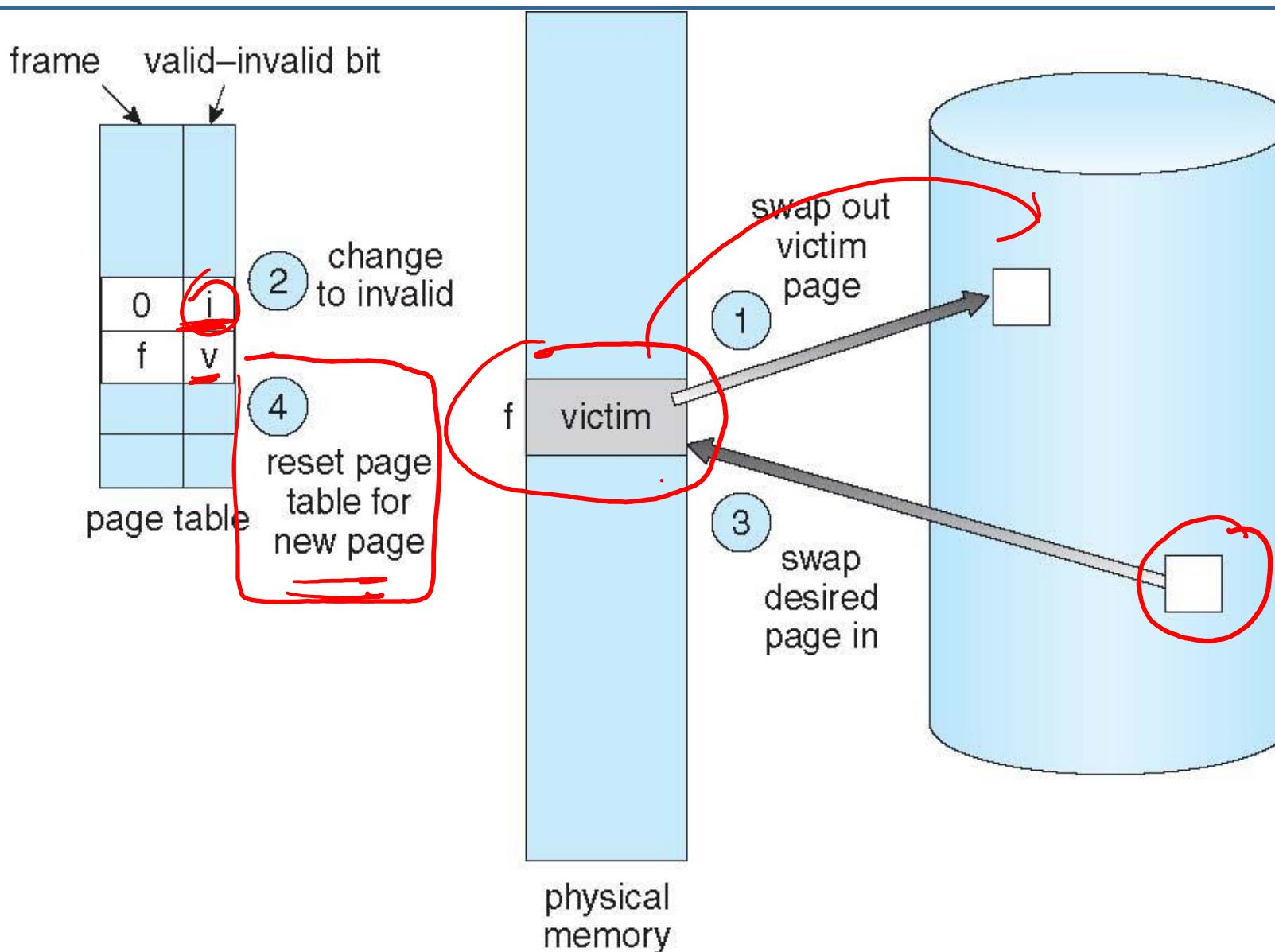
1. Find the location of the desired page on disk
2. Find a free frame:
 - If there is a free frame, use it
 - If there is no free frame, use a page replacement algorithm to select a **victim frame**
 - Write victim frame to disk if dirty = 1 → $\rightarrow \text{W}^-$
3. Bring the desired page into the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

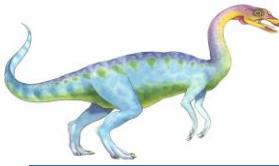
Note now potentially 2 page transfers for page fault – increasing EAT





Page Replacement





PAGE REPLACEMENT ALGORITHMS

Finding the “**victim**” page





→ demand paging

Page and Frame Replacement Algorithms

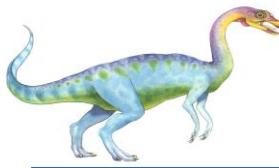
- **Frame-allocation algorithm** determines
 - How many frames to give each process
 - Which frames to replace
- **Page-replacement algorithm**
 - Want lowest page-fault rate on both first access and re-access
 - Evaluate algorithm by running it on a particular string of memory references (reference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

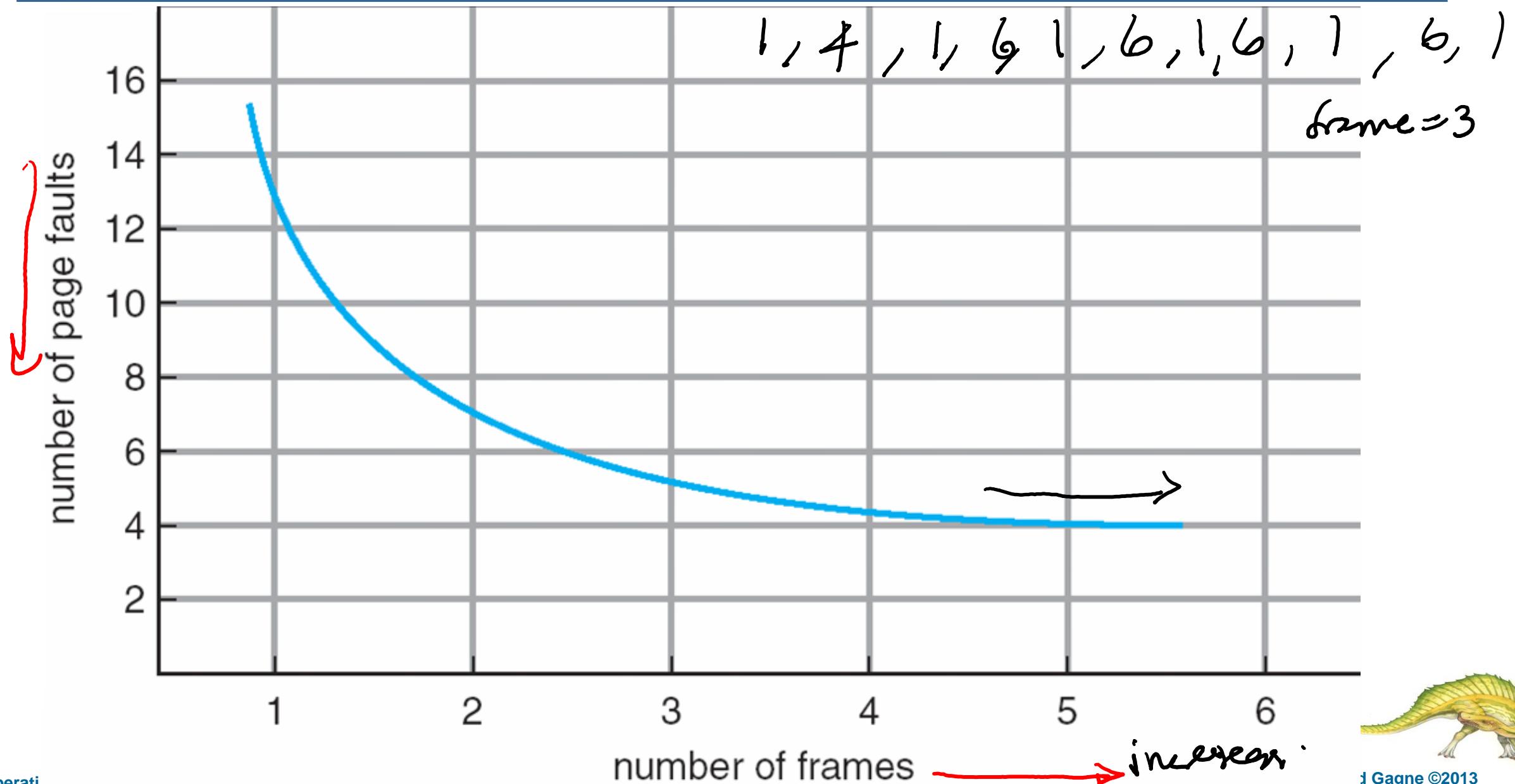
0100 0432 0101
—
0612 0102 0103
—
in bytes 1 page

1 4 1
6 1 1





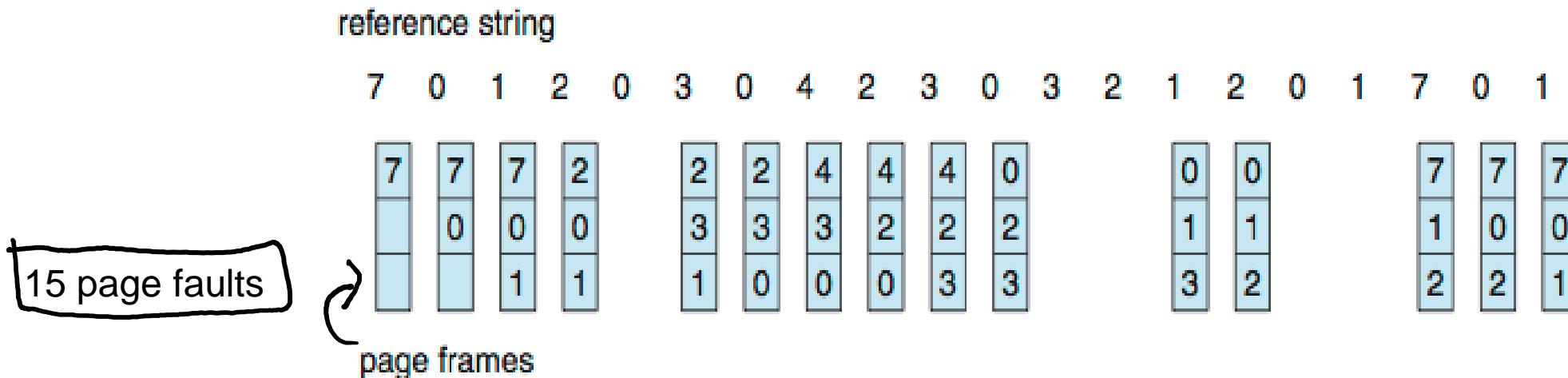
Graph of Page Faults Versus The Number of Frames





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)



- 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!
 - ▶ **Belady's Anomaly**
- How to track ages of pages?
 - Just use a FIFO queue



FIFO

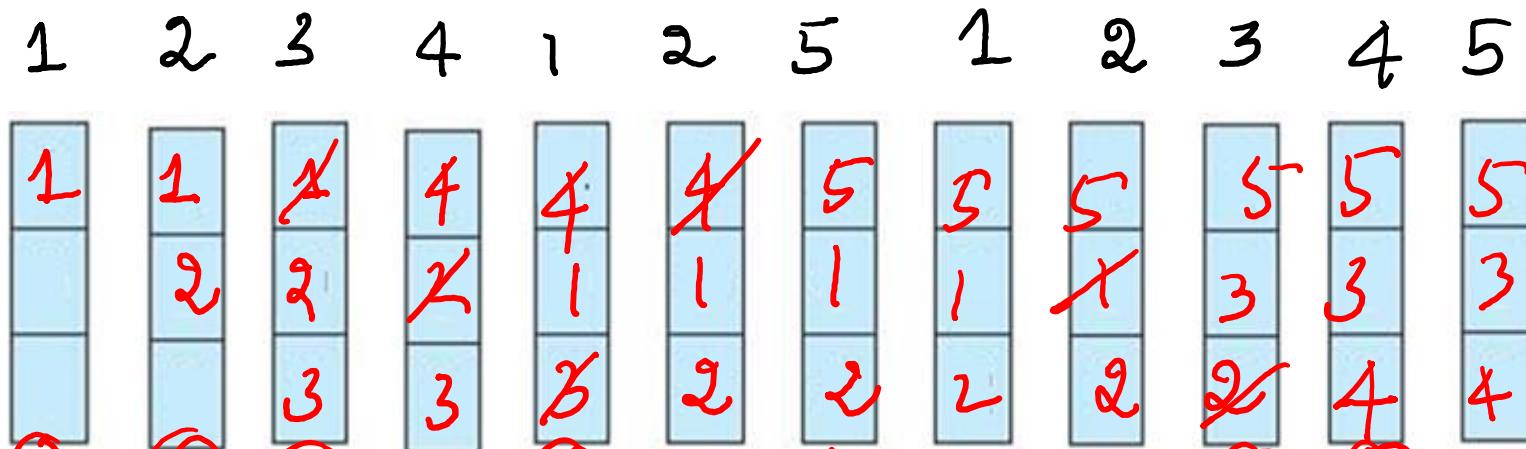
reference string		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	→	7	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	7	0	1
{		7	7	X	2	2	X	4	4	4	0	0	0	0	0	0	0	7	7	7	
x1	→	0	0	→ 1	0	0	3	3	3	2	2	2	2	2	2	2	1	0	0	0	
x2	→	1	4	Hit	5	X	6	7	8	9	10	11	X	11	12	13	14	15	2	1	
x3																					
page frames																					

Frame size = 3

Total No. of page faults = 15

No. of hits = 5 \Rightarrow Hit Ratio = $\frac{5}{20} = \frac{1}{4} = 25\%$

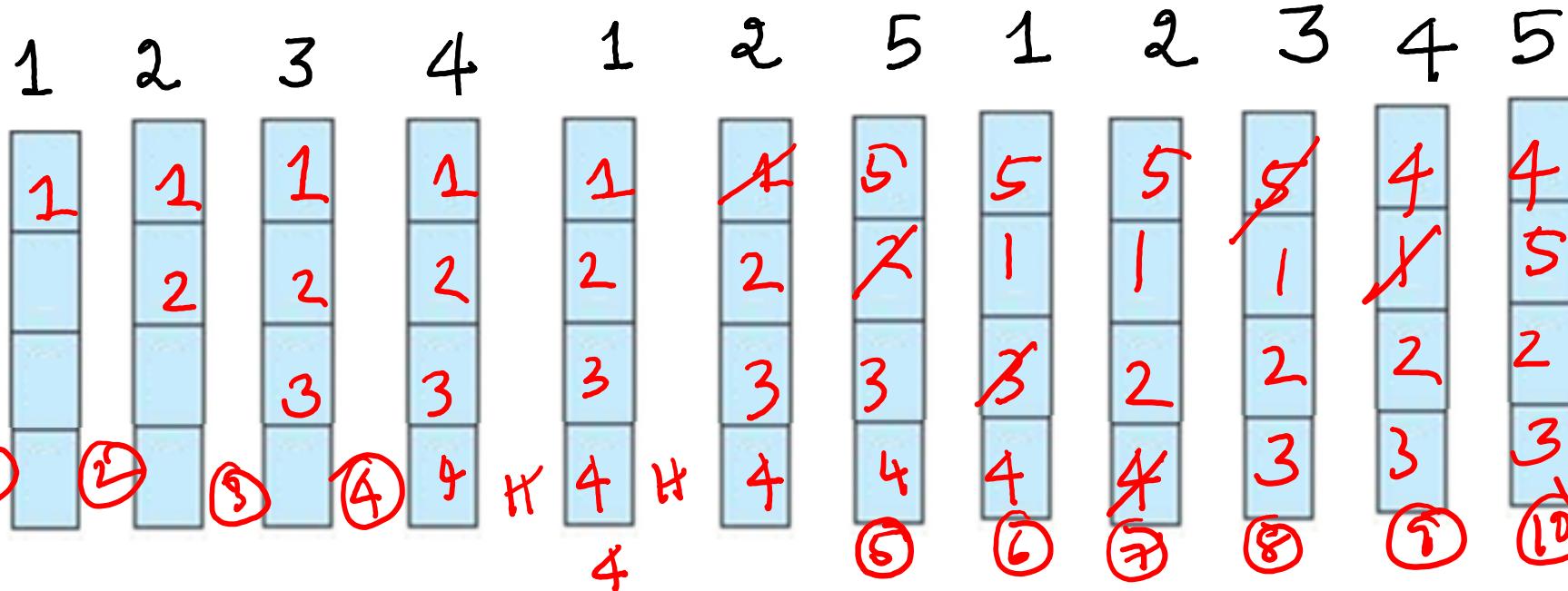
Reference String :



Page faults = ?

9

frame miss = 3



Page faults = ?

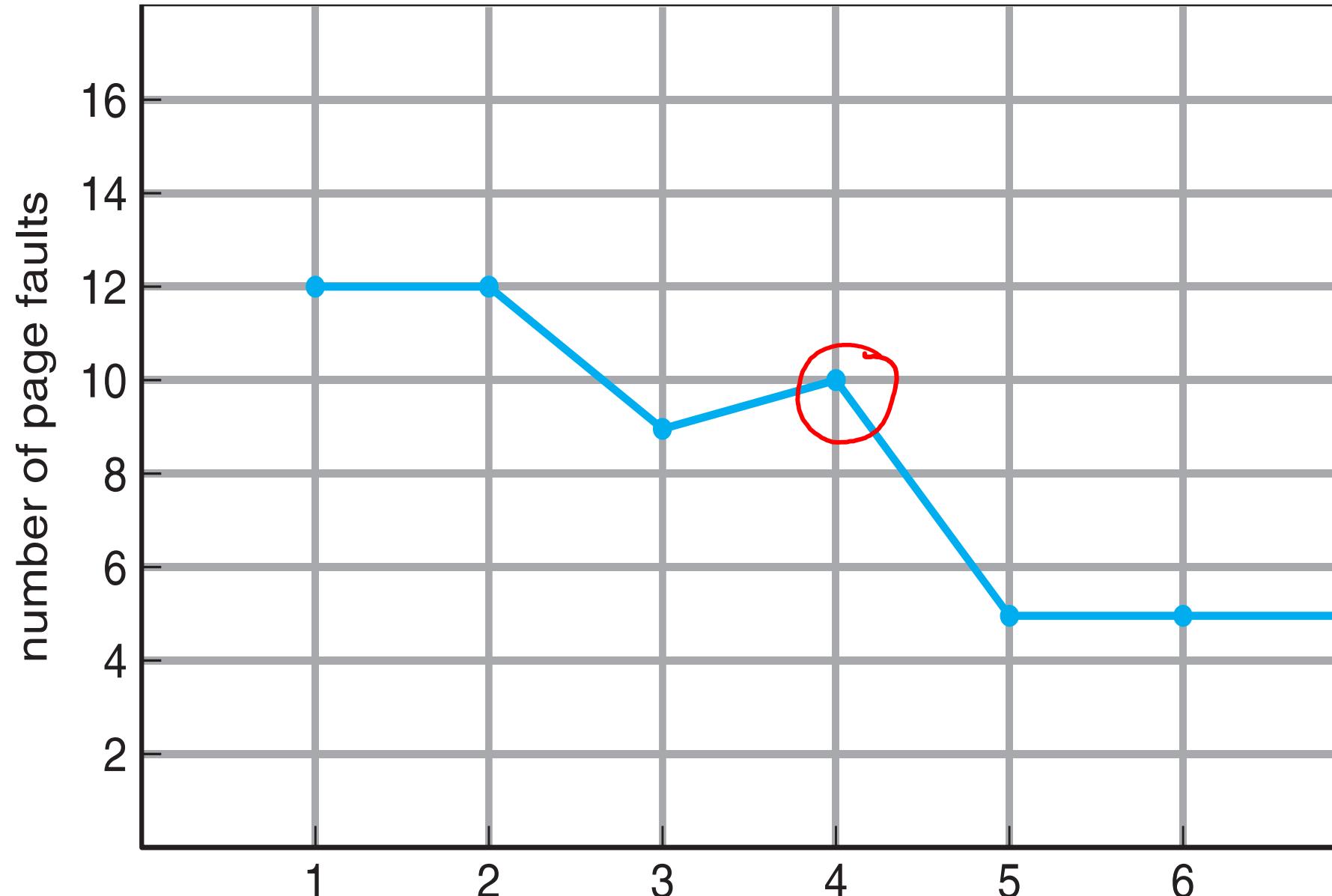
10

$10 > 9$
(4) > (3)

Belady's Anomaly



FIFO Illustrating Belady's Anomaly



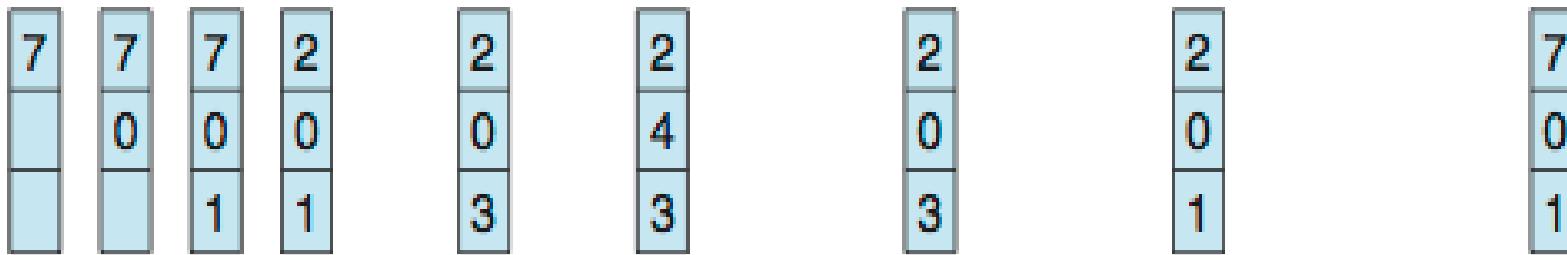


Optimal Algorithm / Page Replacement

- 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

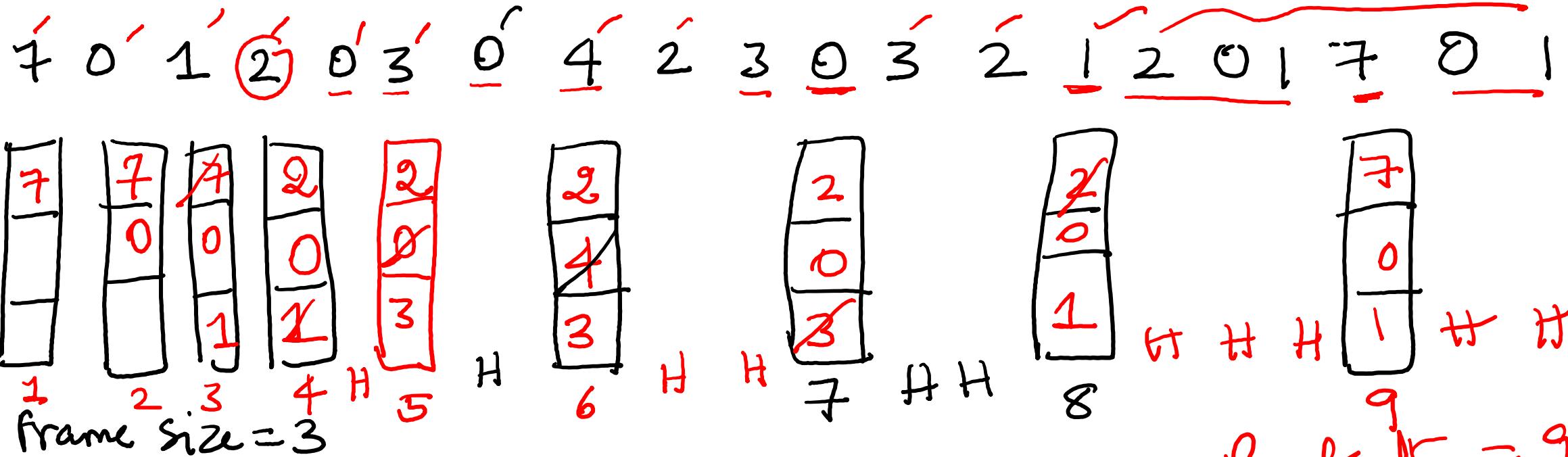
reference string

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

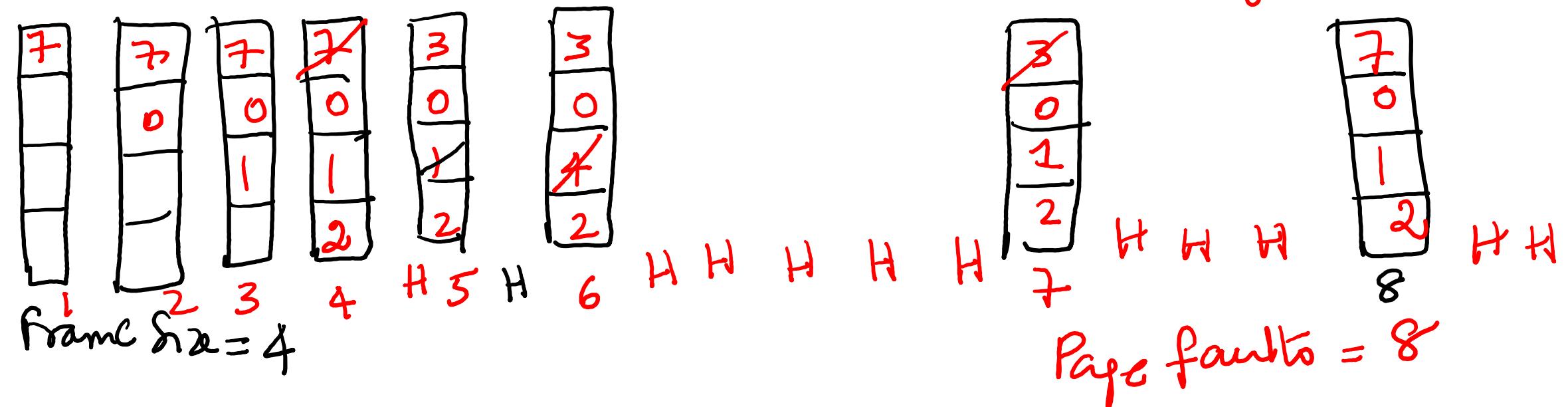


page frames





Page faults = 9 .

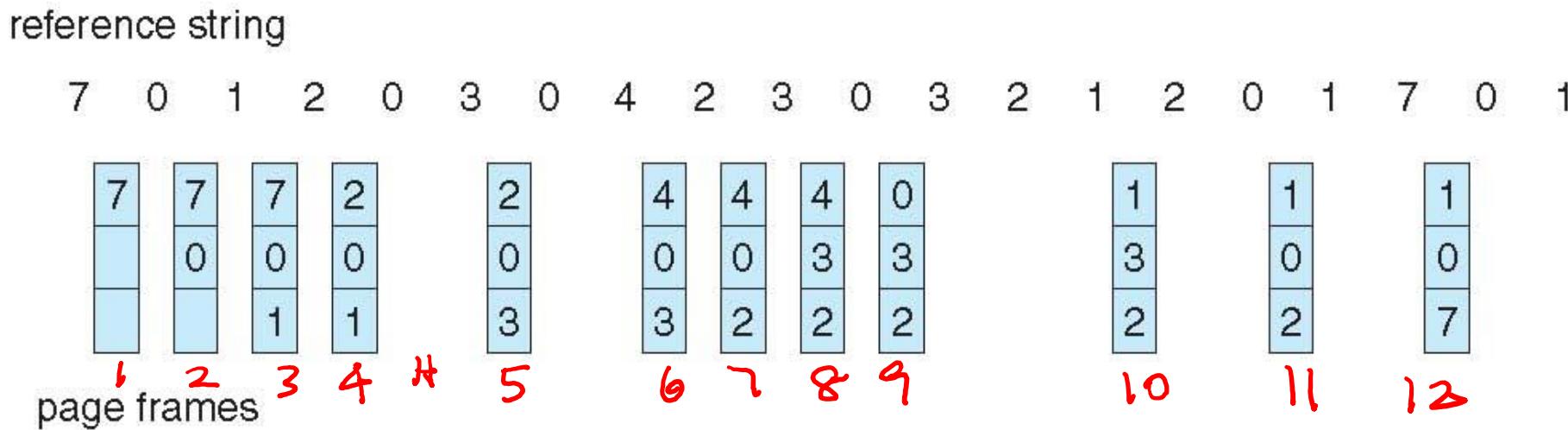


Page faults = 8



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



- 12 faults – better than FIFO but worse than OPT
- Generally good algorithm and frequently used
- But how to implement?



reference string

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
7	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	7	0	1

page frames

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

No. of page frames = 3

Total no. of page faults
= 12

Liu

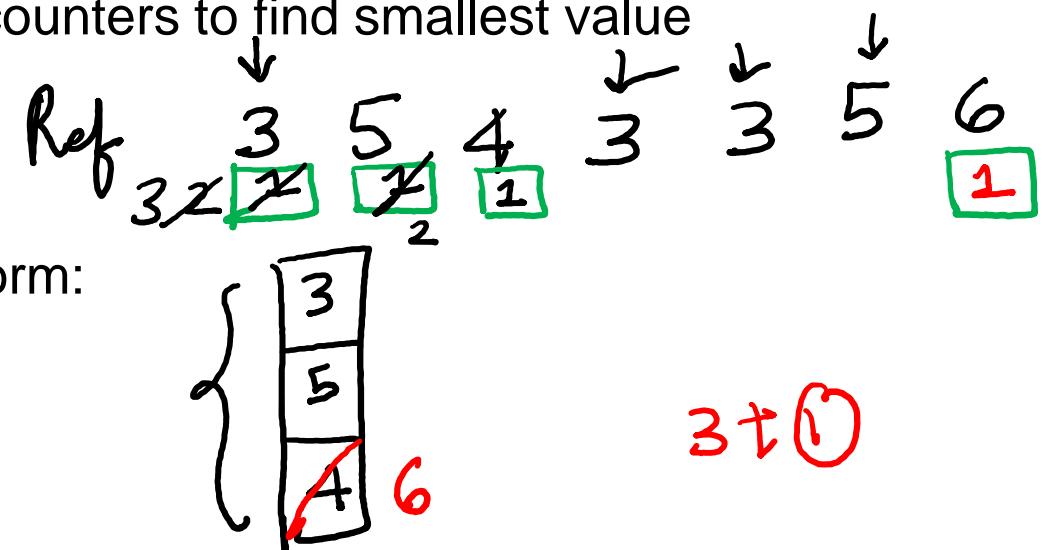
$$9 < 12 < 15$$

OPT FIFO



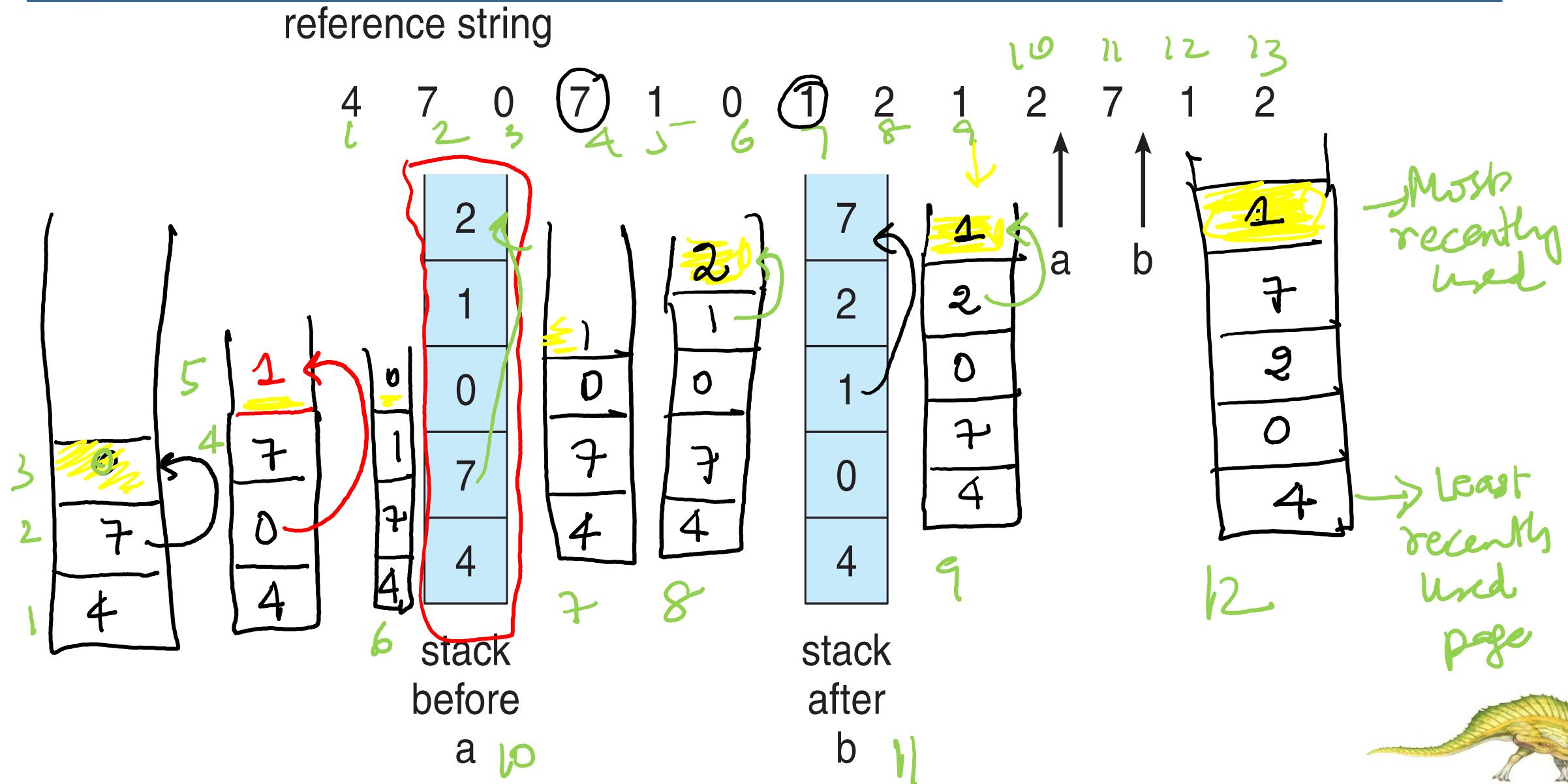
LRU Algorithm (Cont.)

- 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
 - ▶ requires 6 pointers to be changed
 - But each update more expensive
 - No search for replacement
- LRU and OPT are cases of **stack algorithms** that don't have Belady's Anomaly





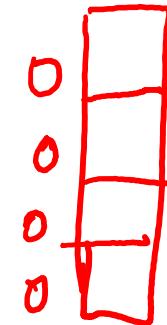
Use Of A Stack to Record Most Recent Page References



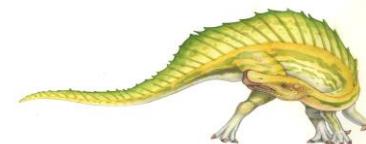


LRU Approximation Algorithms

- 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)
 - ▶ We do not know the order, however
- **Second-chance algorithm** → FIFO
 - 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

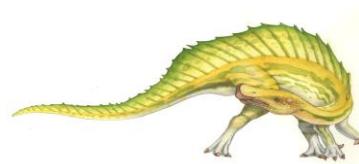
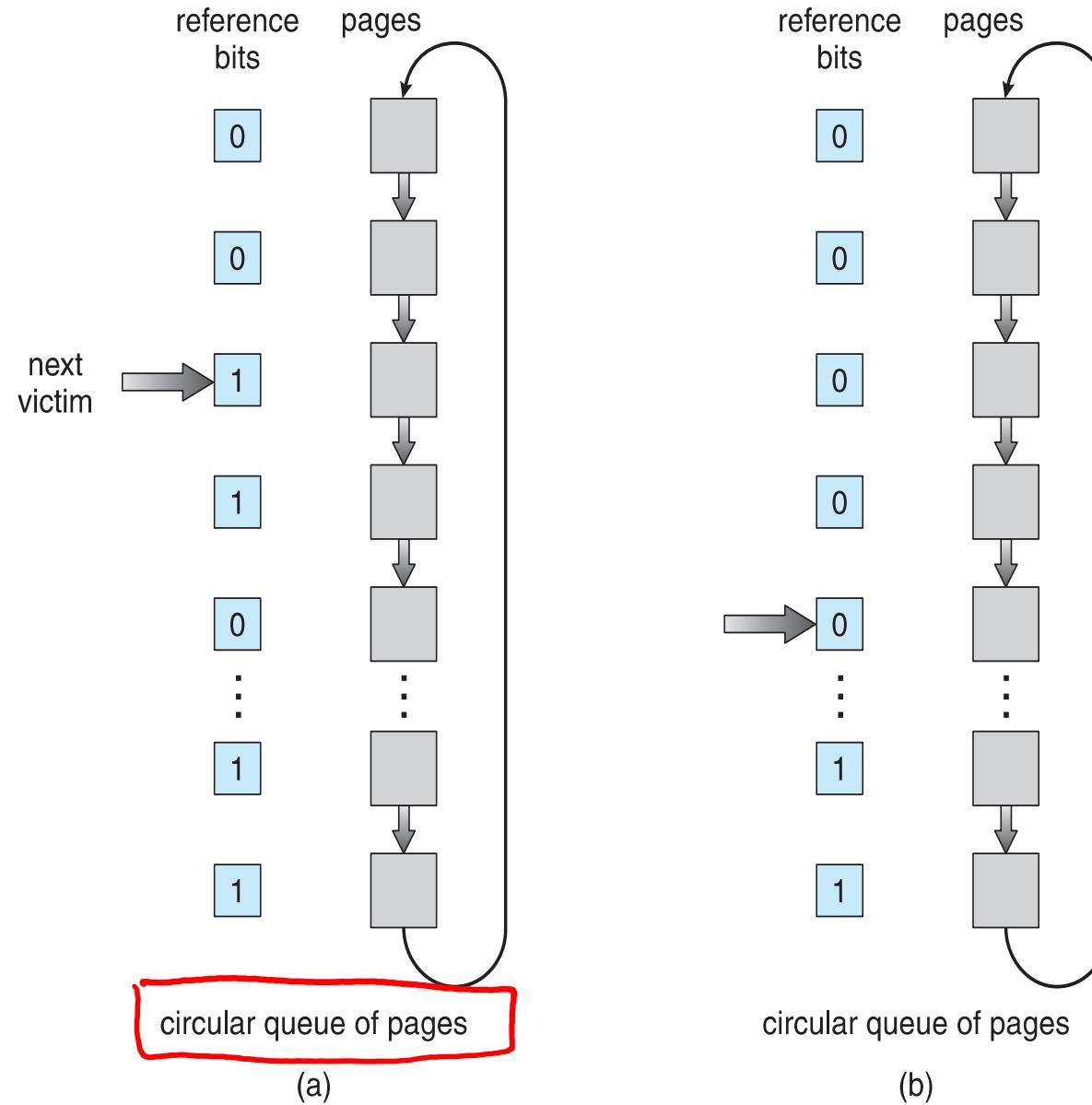


// 2nd chance





Second-Chance (clock) Page-Replacement Algorithm



reference string

ref bit	7	0	1	2	0	3	0	4	2	3	0	3	2	1	2	0	1	7	0	1
0	7	0	7	0	2	0	3	0	2	0	4	0	4	0	4	0	1	1	1	1
0	0	0	0	0	0	1	0	0	1	0	0	3	0	3	1	3	0	0	1	0
0	0	0	1	0	1	0	X	0	3	0	2	0	2	0	2	1	2	0	7	0

1 2 3 4 X 5 X 6 7 8 9 X X 10 X 11 X 12 X X

frame size = 3

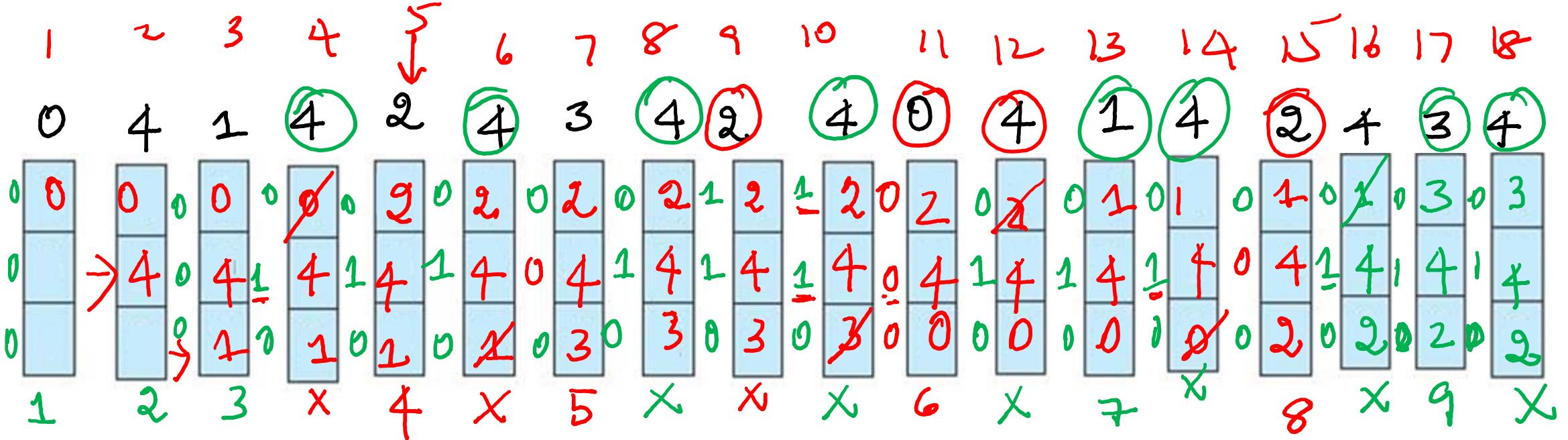
FIFO + ref bit

Page faults = 12

frame size = 4

4	5	5	5	5	6	7	8	6	7	8	6	7	8	4	1	5	2	4	4	1
0	4	0	4	0	4	0	4	0	4	0	8	0	8	1	8	0	8	0	2	2
0	5	1	5	1	5	1	5	1	5	1	5	1	5	1	5	0	4	0	4	1
0	0	0	0	0	0	0	0	0	0	0	6	1	6	1	6	1	6	0	5	0

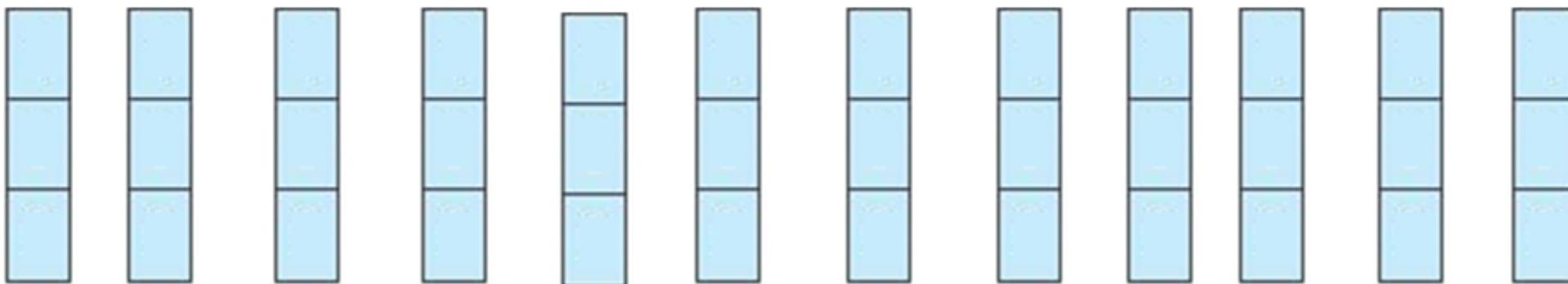
1 2 X 3 X 4 5 X 6 X 7 X 8 X 9 X X X X



Page frames = 3

H/W Problem :

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



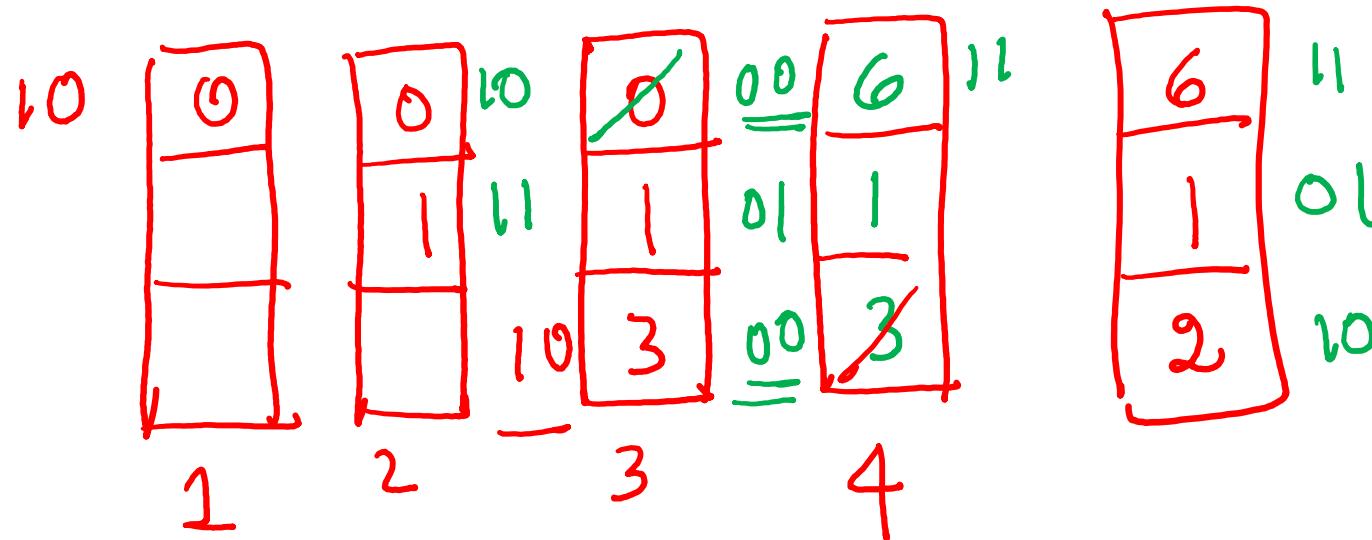


Enhanced Second-Chance Algorithm

- Improve algorithm by using reference bit and modify bit (if available) in concert
- Take ordered pair (reference, modify) *class*
 1. **(0, 0)** neither *recently used not modified* – best page to replace
 2. **(0, 1)** *not recently used but modified* – not quite as good, must write out before replacement *<low priority>*
 3. **(1, 0)** *recently used but clean* – probably will be used again soon
 4. **(1, 1)** *recently used and modified* – probably will be used again soon and need to write out before replacement *2nd chance*
- 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



Ref strip:



frame size = 3

Identify 0 value frame
 $\langle 0, 0 \rangle \rightarrow$ result of the
per step where $\langle 0, 0 \rangle$
not found set bits = 0

$\langle 0, 0 \rangle$
 $\langle 0, 1 \rangle$
 $\langle 1, 0 \rangle$
 $\langle 1, 1 \rangle$



Counting Algorithms

- Keep a counter of the number of references that have been made to each page

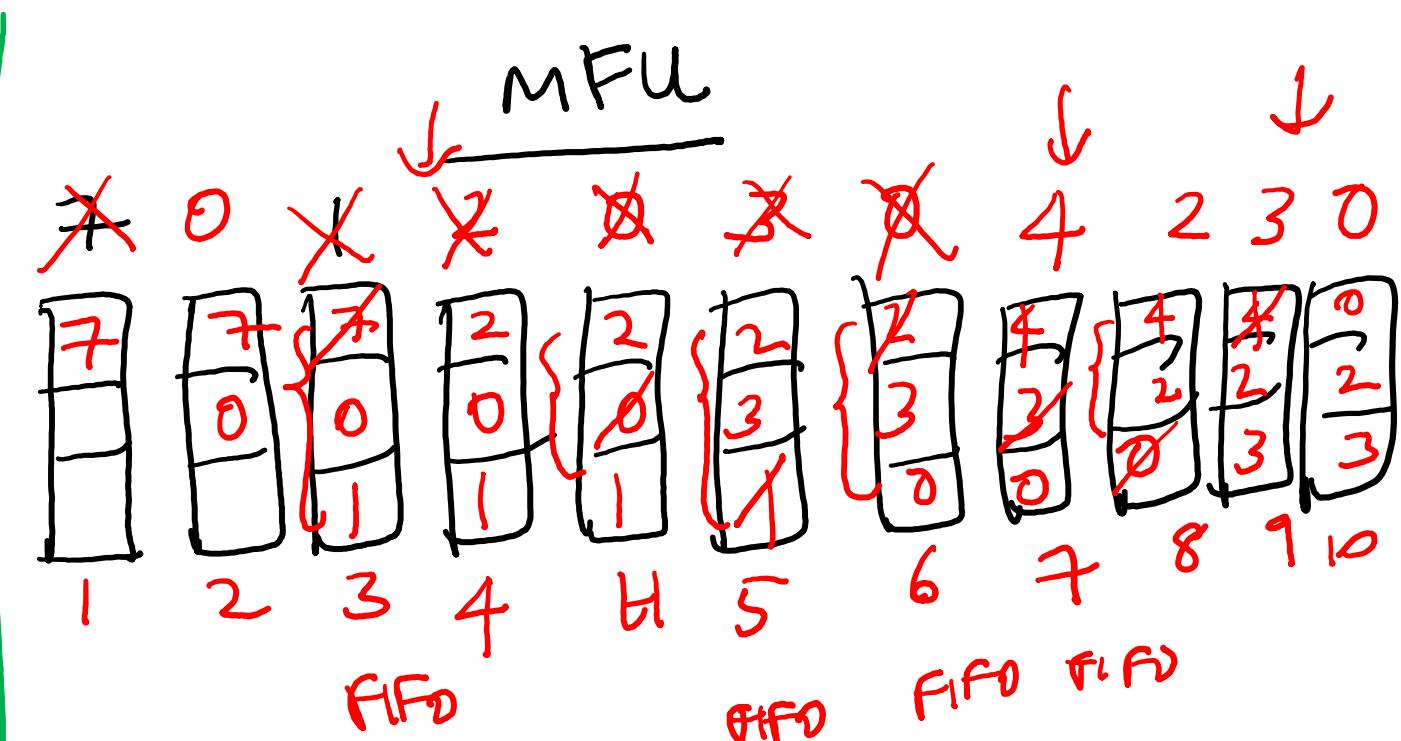
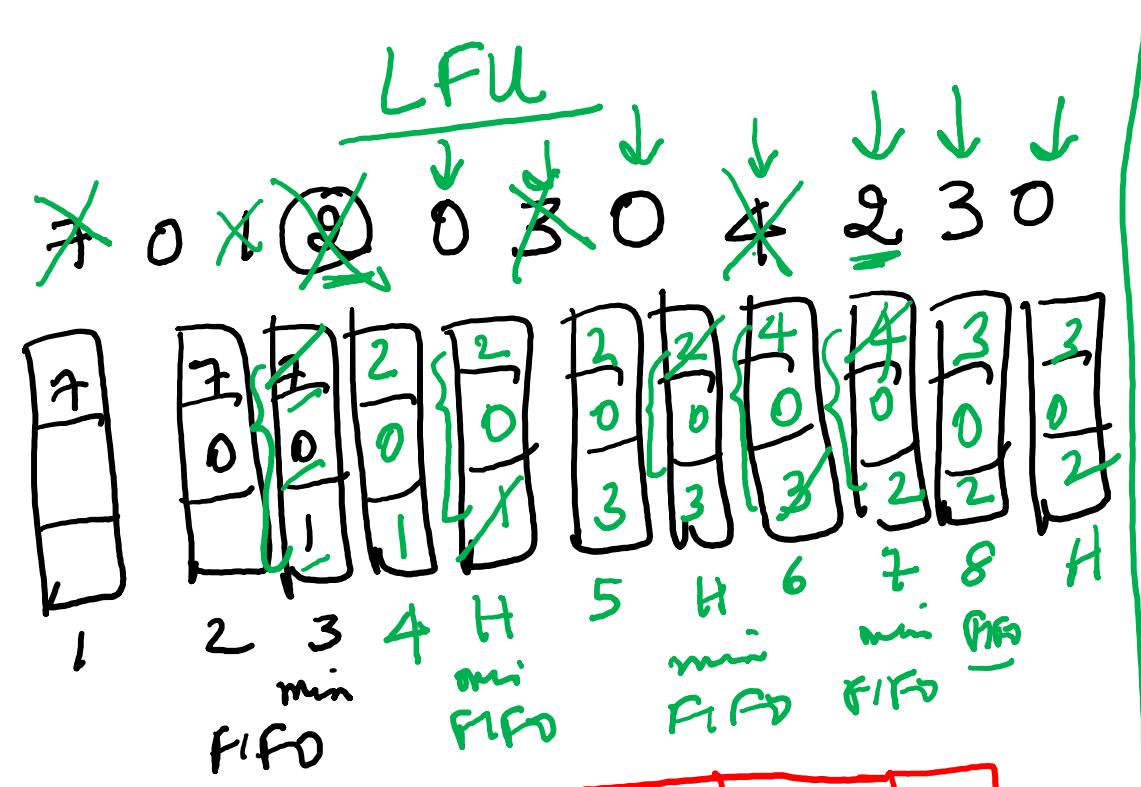
✓ Not common

- **Least Frequently Used (LFU) Algorithm:** replaces page with smallest count

LRU, X

- **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





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

Rules:

① Page smallest freq / count (replace)

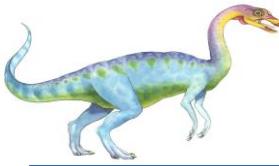
② Page arriving, freq ↑

③ Page leaves, freq ↘
freq tie ⇒ FIFO rule

$$P_{FF} = 8$$

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

1 2 3 4 H 5 6 7 8 9 10



Page-Buffering Algorithms

- Keep a pool of free frames, always
 - Then frame available when needed, not found at fault time
 - Read page into free frame and select victim to evict and add to free pool
 - When convenient, evict victim
- Possibly, keep list of modified pages
 - When backing store otherwise idle, write pages there and set to non-dirty I/O
- 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

$$m = 93 \quad n = 5 \quad 93/5 = 18 \text{ frames}$$

- ① □ 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 $m/n \text{ frames}$

- ② □ Proportional allocation – Allocate according to the size of process
- Dynamic as degree of multiprogramming, process sizes change

m frames
 n processes

s_i = size of process p_i

$$S = \sum s_i \quad \frac{10}{137} \quad \frac{127}{137}$$

m = total number of frames

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

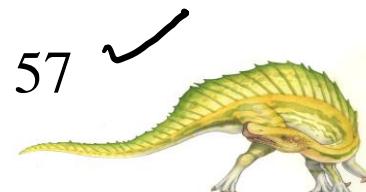
$$m = 64$$

$$s_1 = 10$$

$$s_2 = 127$$

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

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





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
 - 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

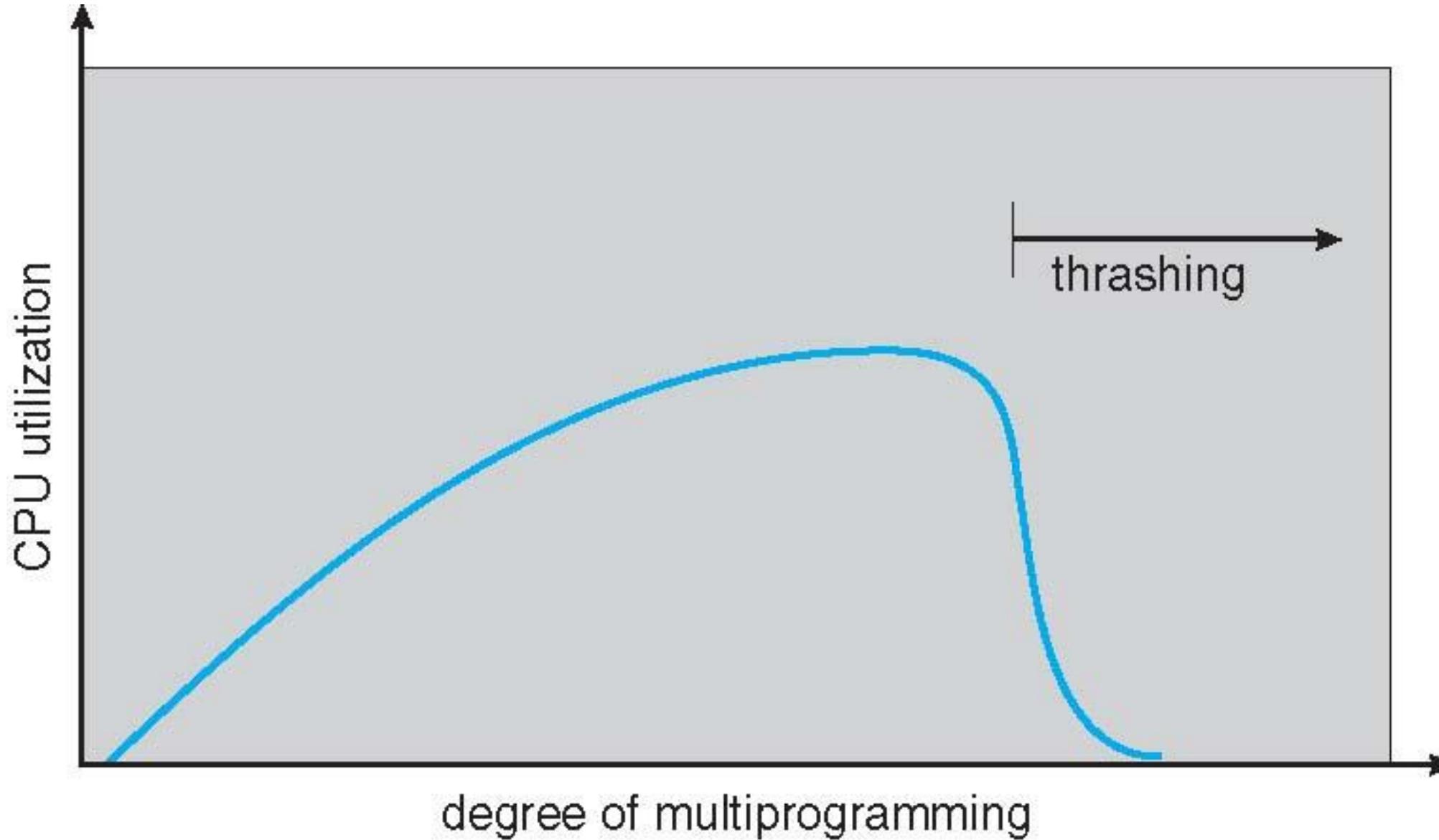
↑ Paging ↓ Executing

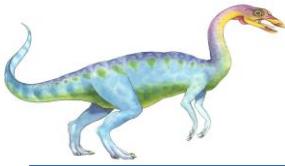
- If a process does not have “enough” pages, the page-fault rate is very high
 - Page fault to get page
 - Replace existing frame
 - But quickly need replaced frame back
 - This leads to:
 - ▶ Low CPU utilization
 - ▶ Operating system thinking that it needs to increase the degree of multiprogramming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





Thrashing (Cont.)





Demand Paging and Thrashing

- Why does demand paging work?

Locality model

set a page that are actively used together

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

- Why does thrashing occur?

$$\Sigma \text{ size of locality} > \text{total memory size}$$

- Limit effects by using local or priority page replacement





Working-Set Model

- ◻ $\Delta \equiv$ working-set window \equiv a fixed number of page references

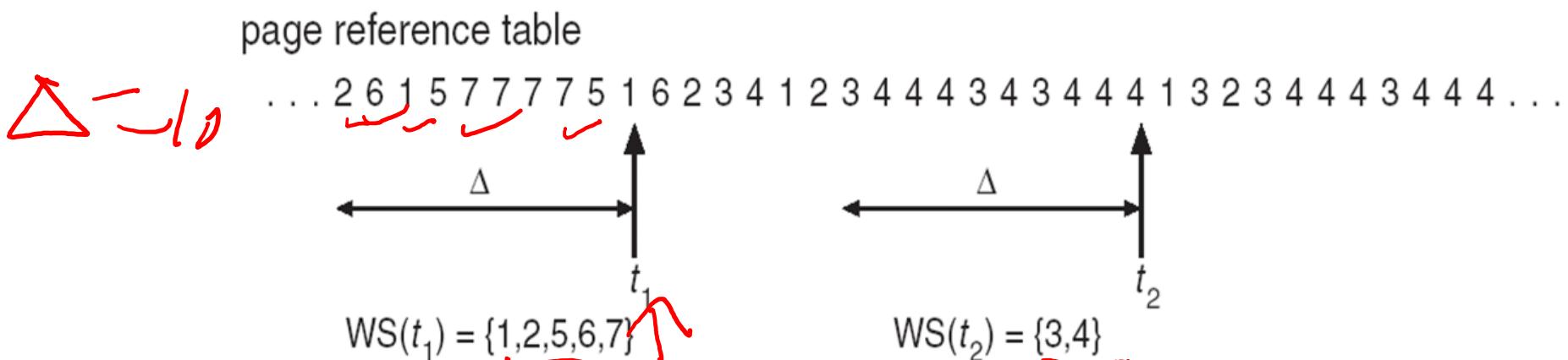
Example: 10,000 instructions

- ◻ WSS_i (working set of Process P_i) =
total number of pages referenced in the most recent 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



if $D > m \Rightarrow$ Thrashing
Policy if $D > m$, then
suspend or swap out one
of the processes





Keeping Track of the Working Set

- Approximate with interval timer + a reference bit
- Example: $\Delta = 10,000$ *Refers*
 - 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

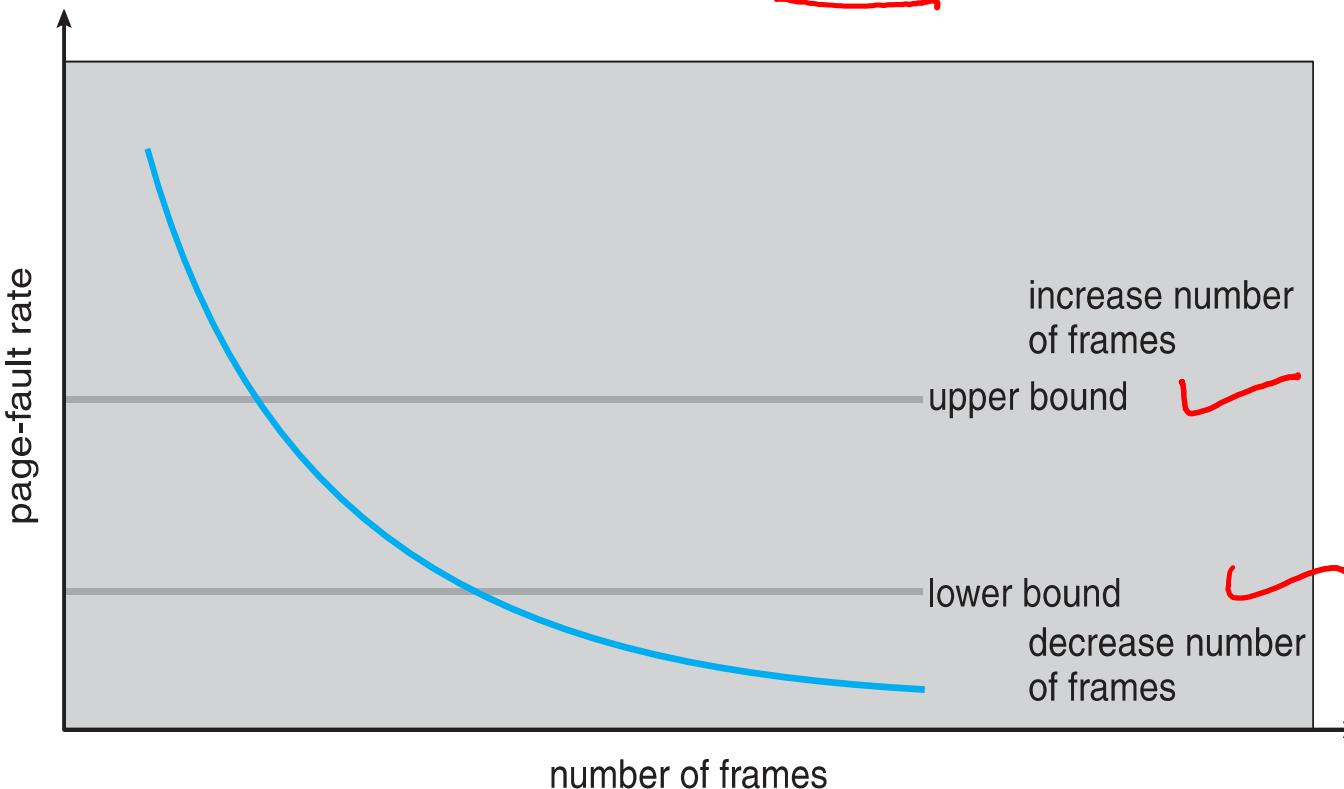




② Page-Fault Frequency

Thresh
PF rate ↑

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



End of Chapter 9

