

CHAPTER 10 - VIRTUAL- MEMORY

OBJECTIVES


- Describe the benefits of a virtual memory system
- Explain the concepts of demand paging, page-replacement algorithms, and allocation of page frames
- Discuss the principle of the working-set model

BACKGROUND

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
 - Program and programs could be larger than physical memory

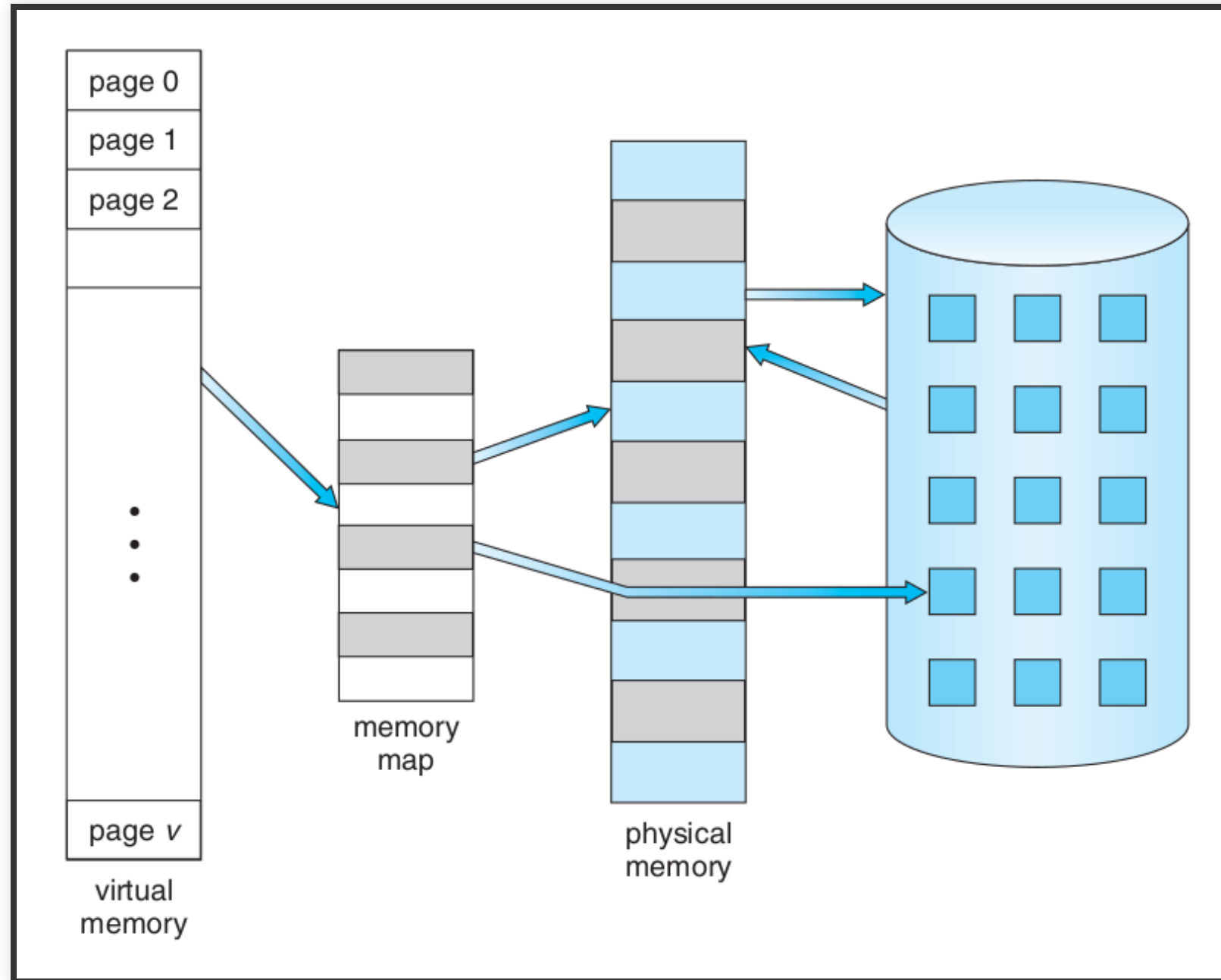
BACKGROUND

 **Virtual memory** – separation of user logical memory from physical memory

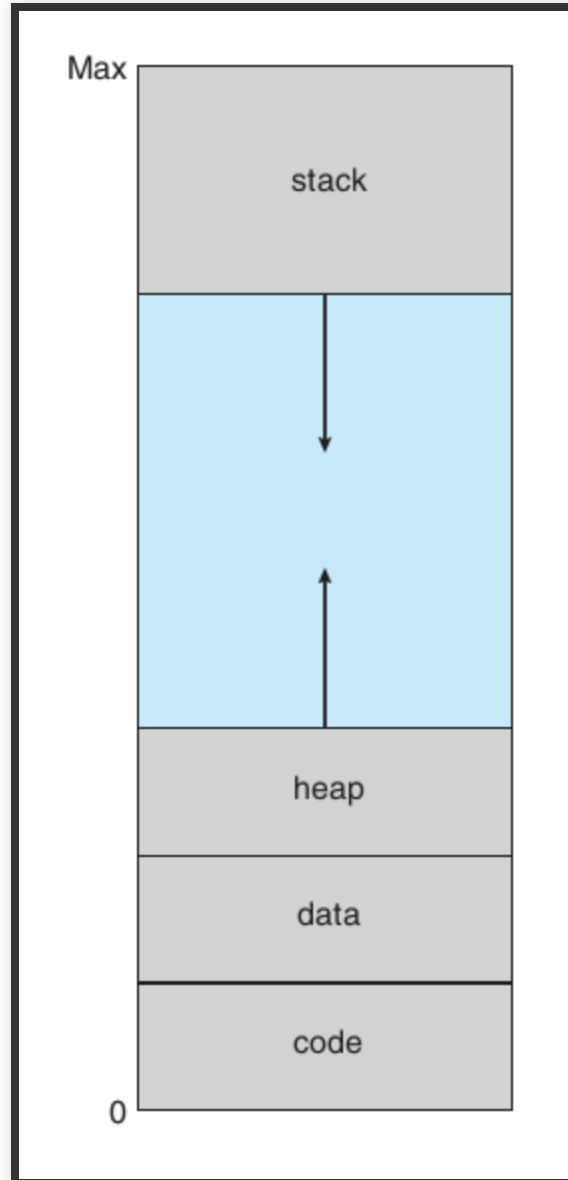
BACKGROUND

- 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

VIRTUAL MEMORY LARGER THAN PHYSICAL MEMORY



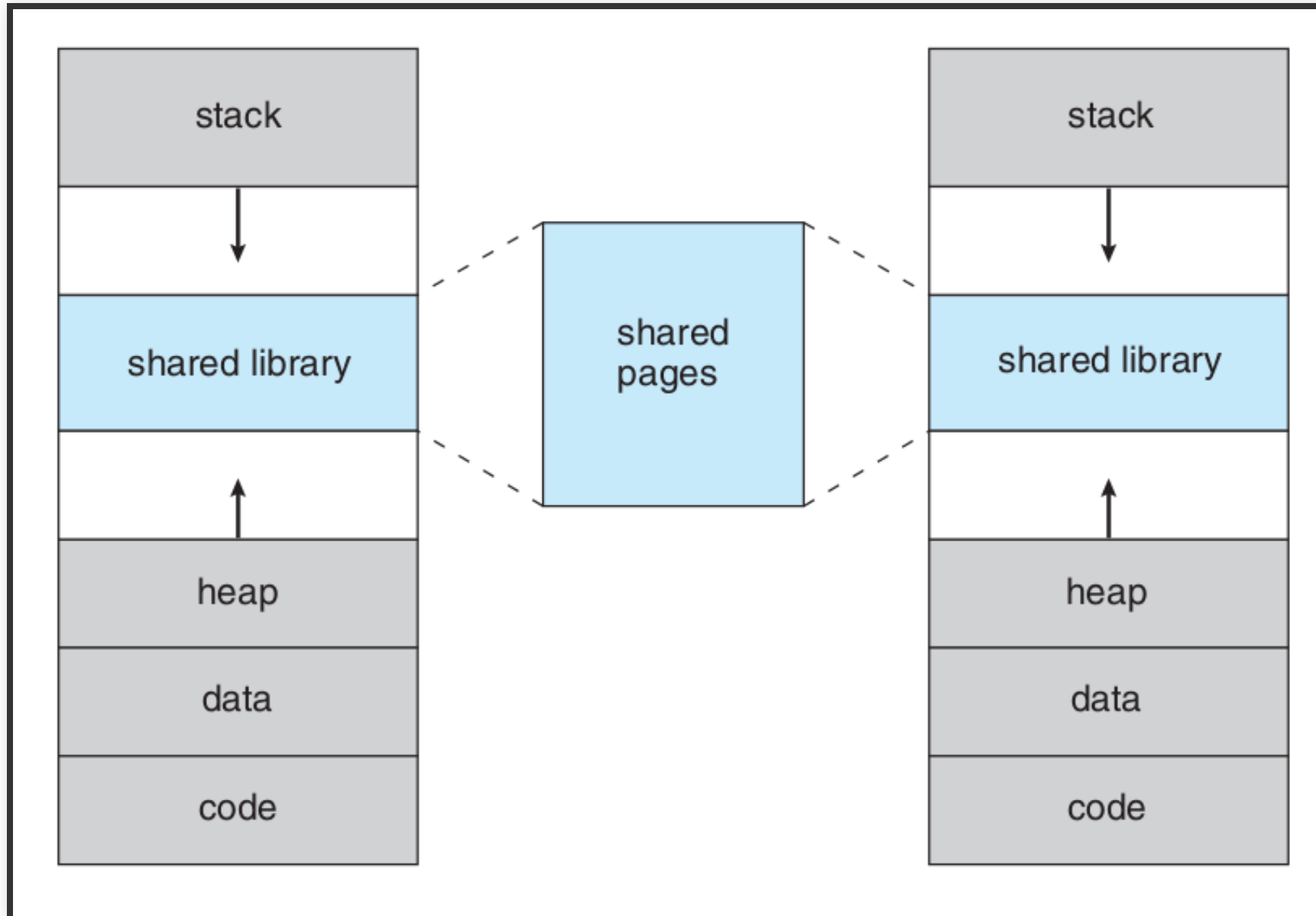
VIRTUAL-ADDRESS SPACE



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 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

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

DEMAND PAGING

- 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 BIT

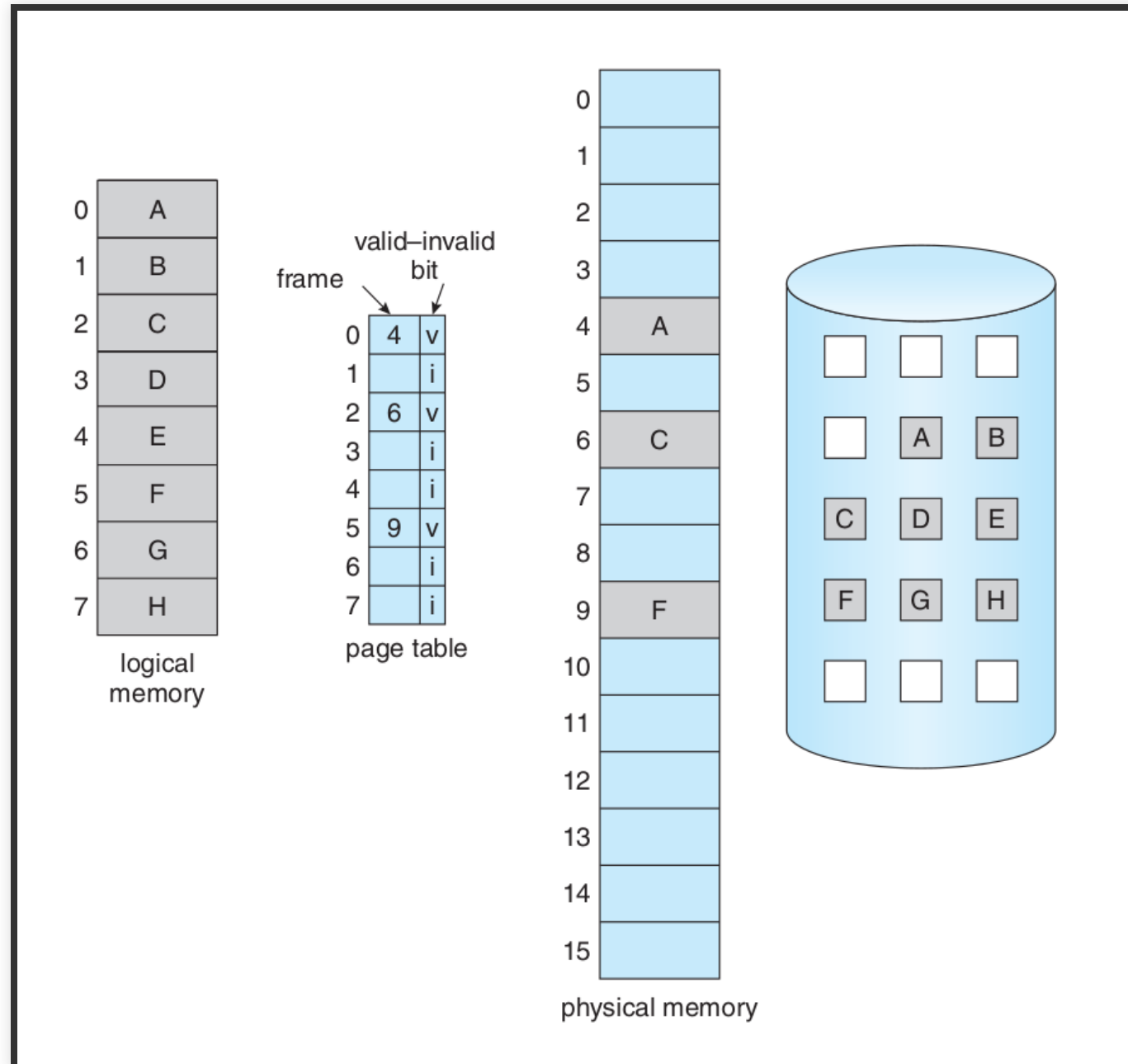
- With each page table entry a valid–invalid bit is associated
 - ($v \Rightarrow$ in-memory, $i \Rightarrow$ not-in-memory)

The diagram shows a table representing a page table. It has two columns: 'Frame #' and 'valid-invalid bit'. The 'valid-invalid bit' column contains a sequence of bits: 'v', 'v', 'v', 'v', 'i', '...', 'i', 'i'. The 'Frame #' column is empty. The table is labeled 'page table' at the bottom.

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

- Initially valid–invalid bit is set to i on all entries
- During address translation, if valid–invalid bit in page table entry is $i \Rightarrow$ **page fault**

PAGE TABLE



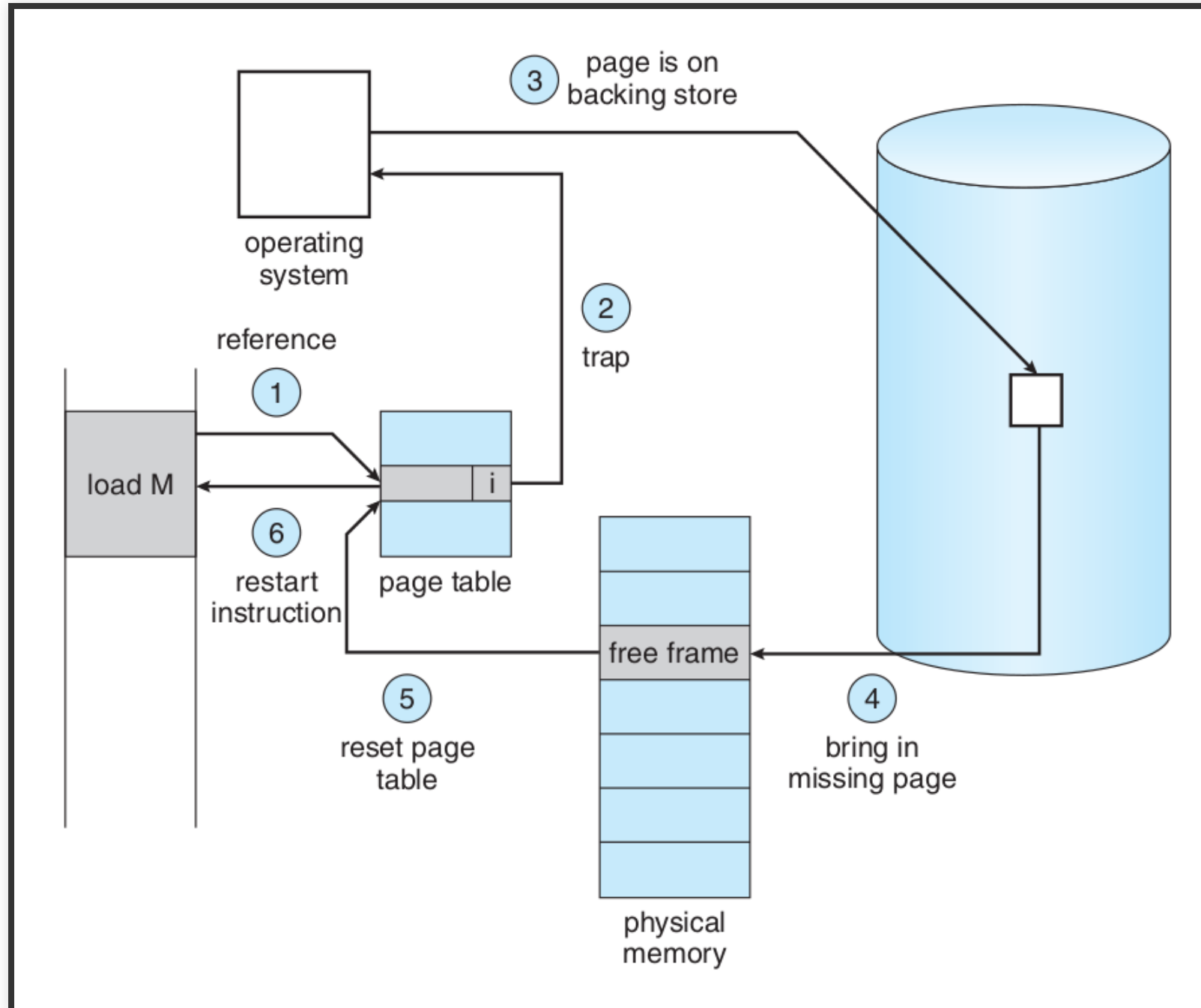
PAGE FAULT

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

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 \rightarrow Set validation bit =
v
5. Restart the instruction that caused the page fault

STEPS IN HANDLING A PAGE FAULT



ASPECTS OF DEMAND PAGING

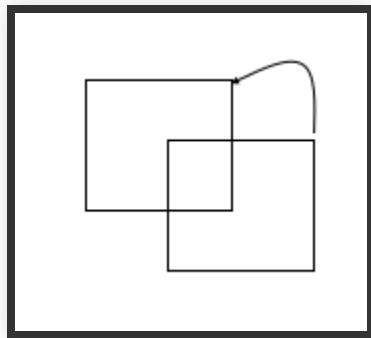
- Extreme case – start process with no pages in memory
 - OS sets instruction pointer to first instruction of process, non-memory-resident → page fault
 - And for every other process pages on first access
 - Pure demand paging
- Actually, a given instruction could access multiple pages → multiple page faults
 - Pain decreased because of locality of reference

ASPECTS OF DEMAND PAGING

- Hardware support needed for demand paging
 - Page table with valid / invalid bit
 - Secondary memory (swap device with swap space)
 - Instruction restart

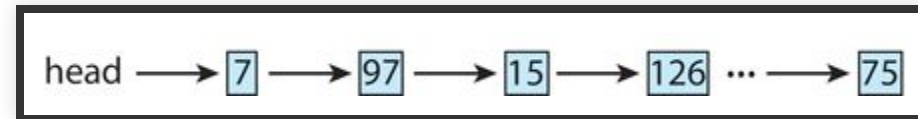
INSTRUCTION RESTART

- Consider an instruction that could access several different locations



- block move
- auto increment/decrement location
- Restart the whole operation?
 - What if source and destination overlap?

FREE-FRAME LIST



- In general, free pages are allocated from a pool of zero-fill-on-demand pages
 - Why zero-out a page before allocating it?

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

PERFORMANCE OF DEMAND PAGING

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

PERFORMANCE OF DEMAND PAGING

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**)

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

DEMAND PAGING EXAMPLE

- Memory access time = 200 nanoseconds
- Average page-fault service time = 8 milliseconds
- $EAT = (1 - p) * 200 + p (8 \text{ milliseconds})$
 $= (1 - p * 200 + p * 8,000,000$
 $= 200 + p * 7,999,800$
- 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!!

DEMAND PAGING EXAMPLE

If want performance degradation < 10 percent

- $220 > 200 + 7,999,800 * p$
 $20 > 7,999,800 * p$
 $p < .0000025$
- < one page fault in every 400,000 memory accesses

DEMAND PAGING OPTIMIZATIONS

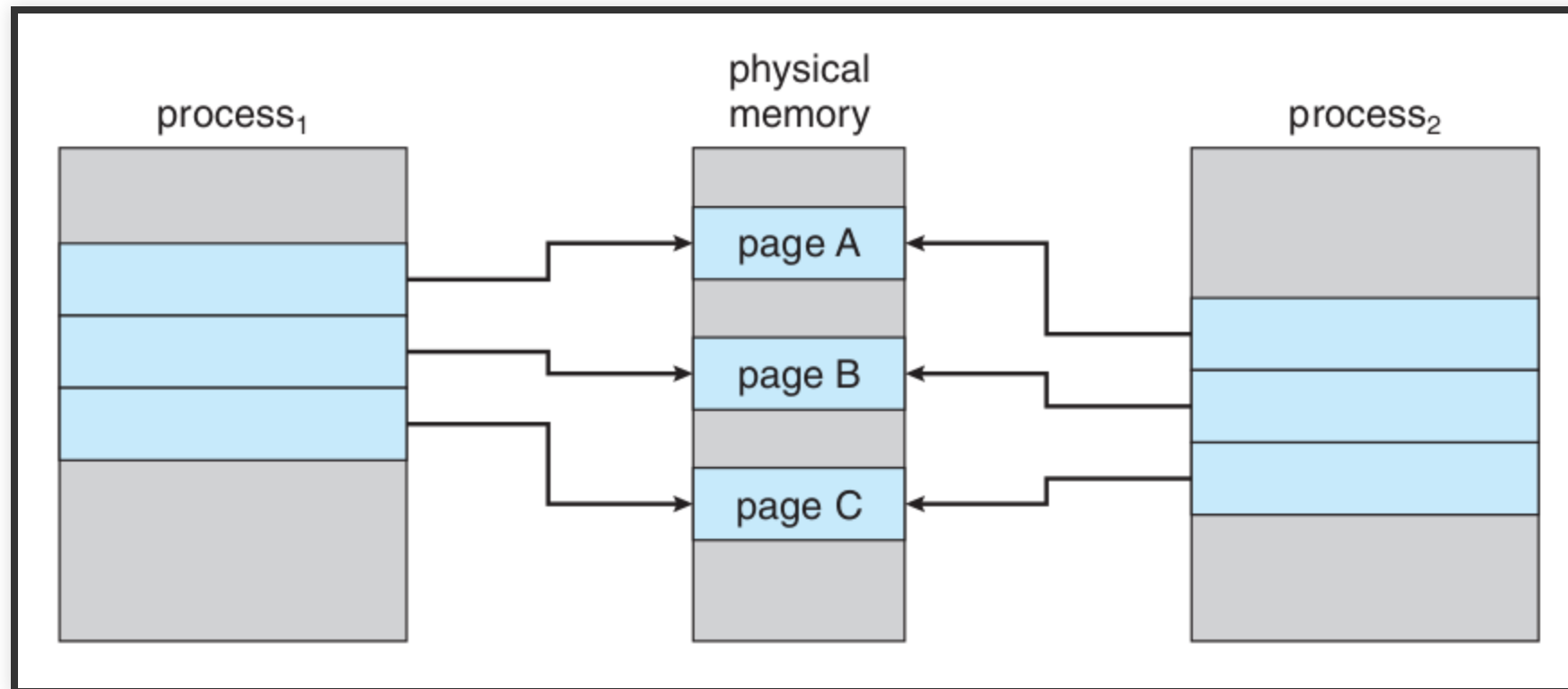
- Copy entire process image to swap space at process load time
 - Then page in and out of swap space
 - Used in older BSD Unix
- Demand page in from program binary on disk, but discard rather than paging out when freeing frame
 - Used in Solaris and current BSD

COPY-ON-WRITE

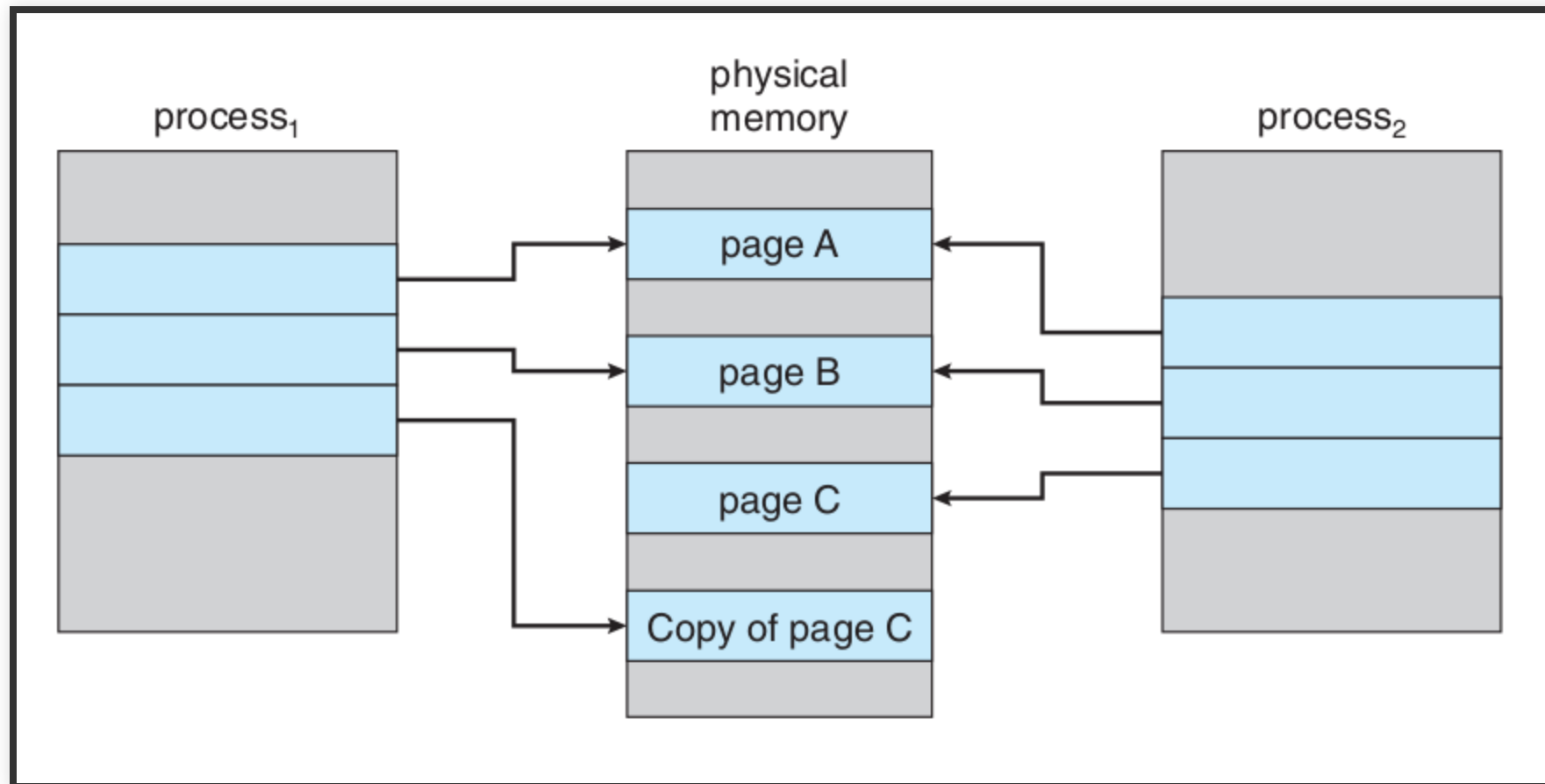
COPY-ON-WRITE

- 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

BEFORE P1 MODIFIES PAGE C



AFTER P1 MODIFIES PAGE C



COPY-ON-WRITE

- `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

PAGE REPLACEMENT

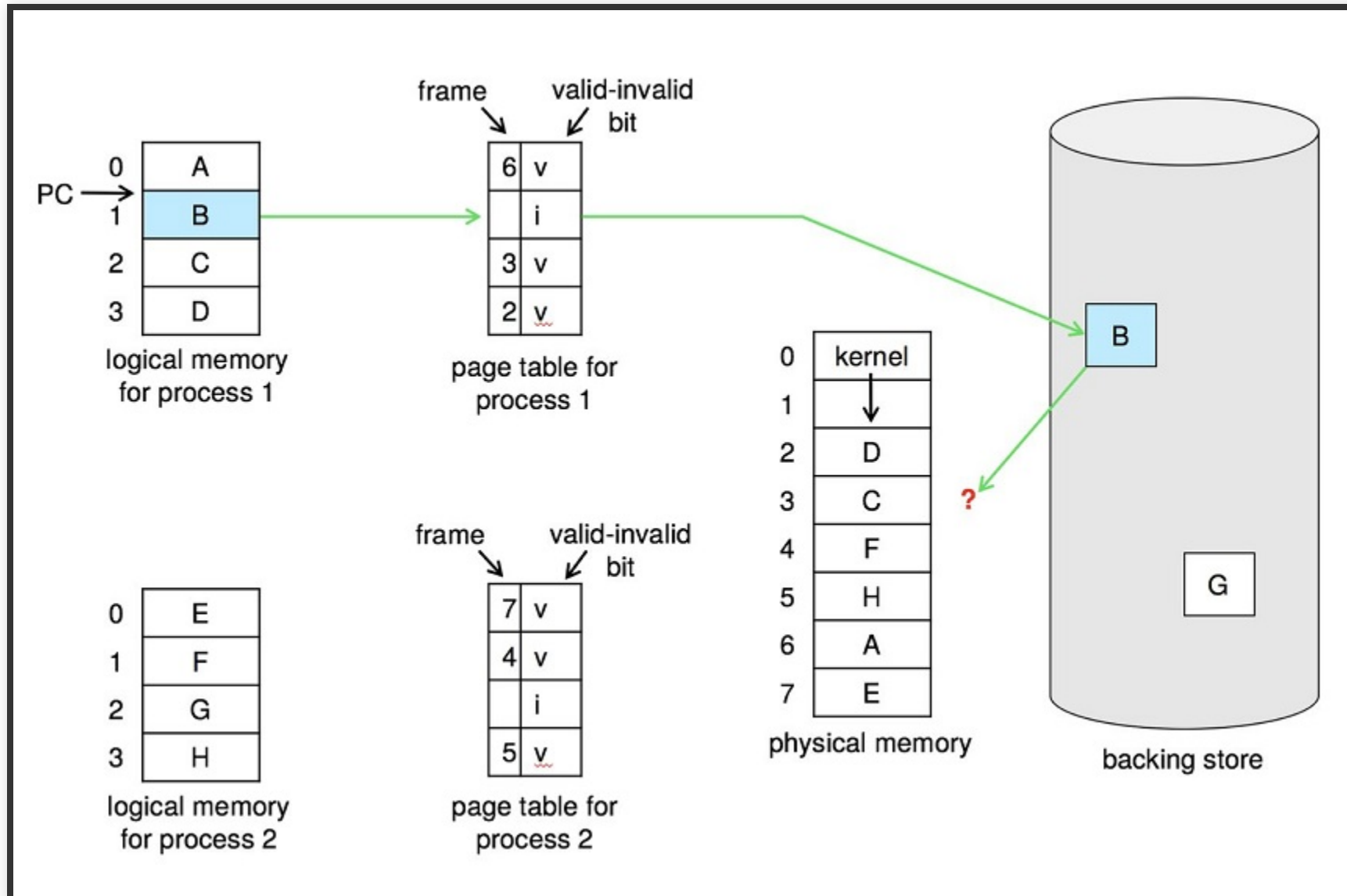
NO FREE FRAME?

- 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
- 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



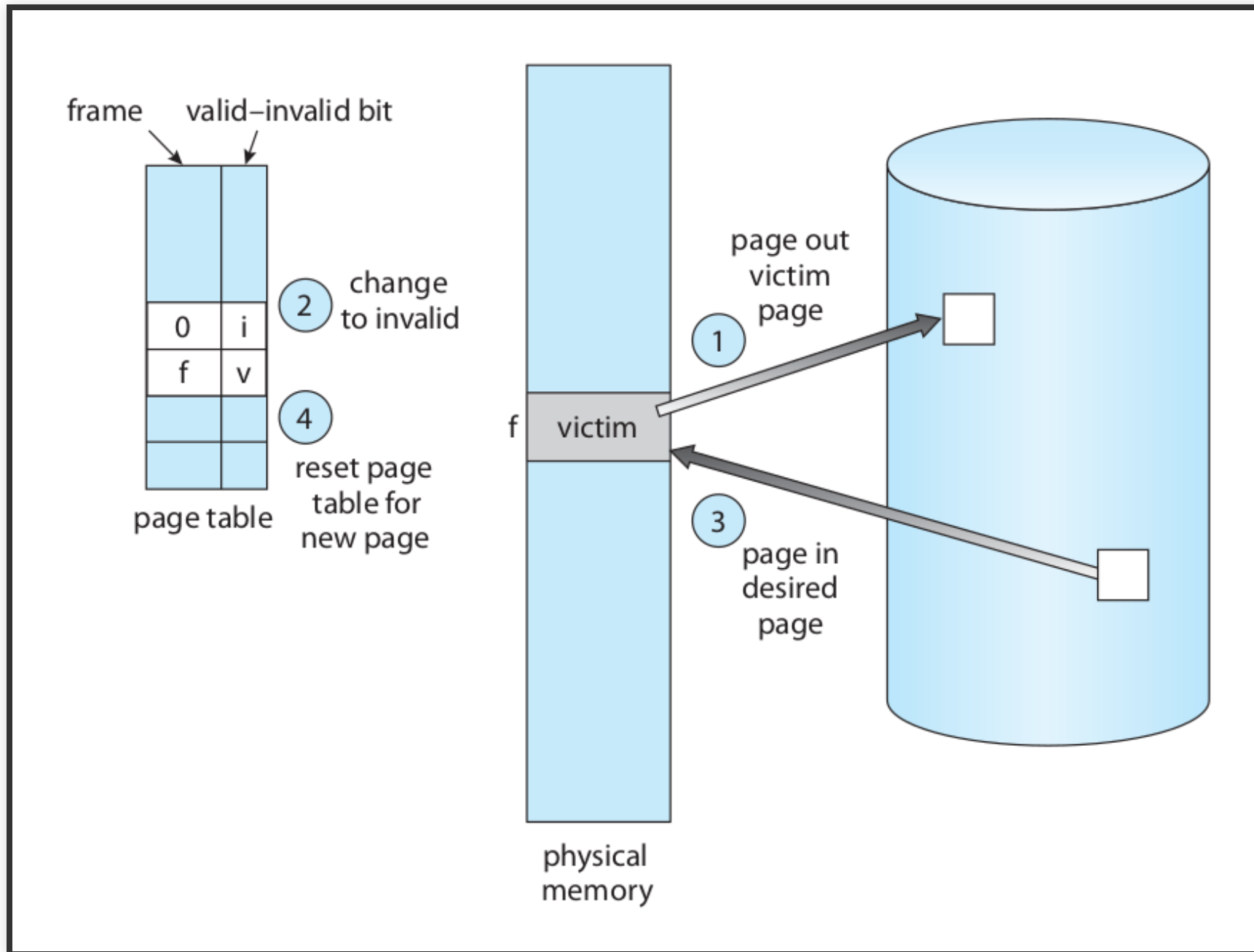
BASIC PAGE REPLACEMENT

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

BASIC PAGE REPLACEMENT

- ⚠ 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

PAGE AND FRAME REPLACEMENT ALGORITHMS

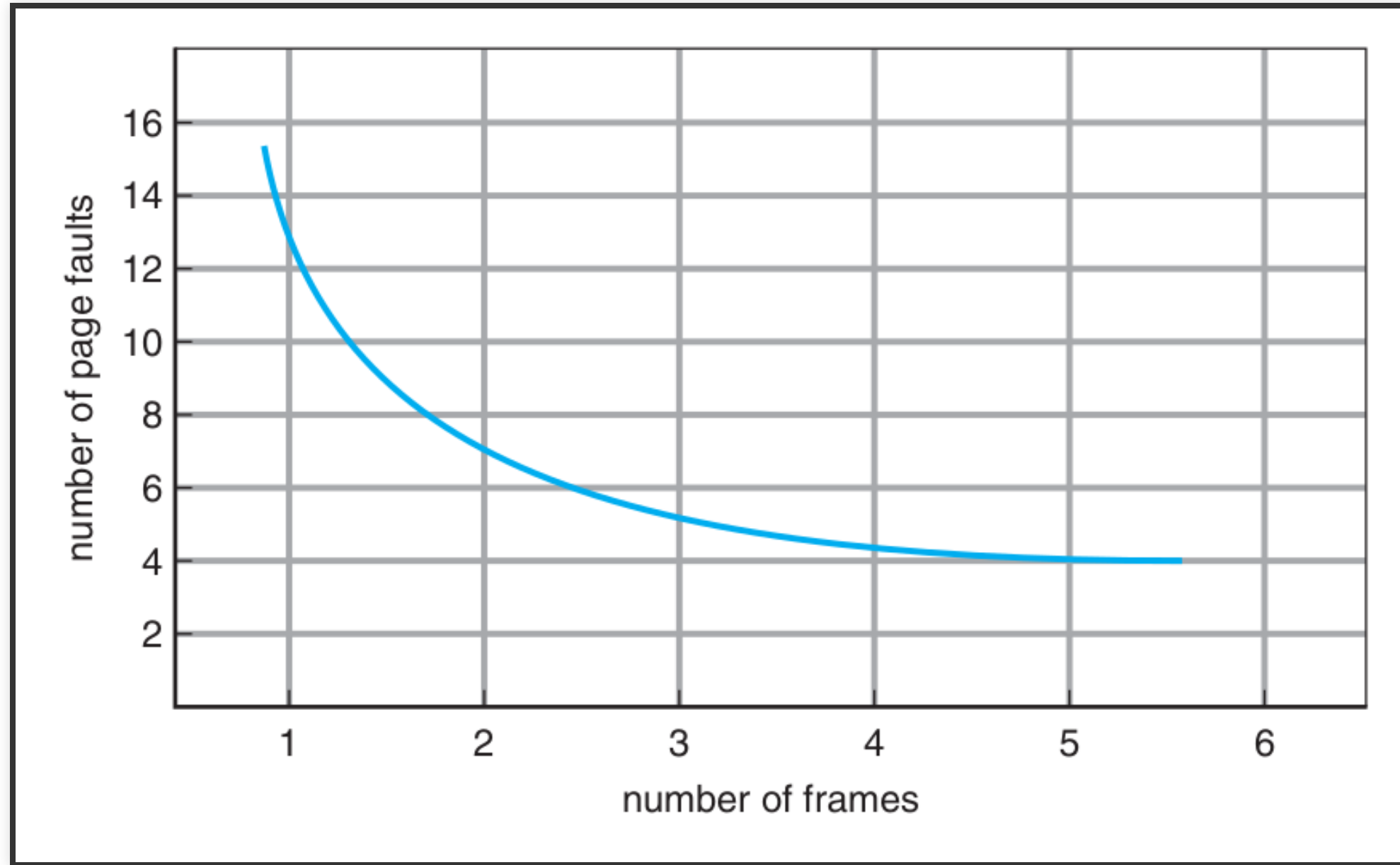
- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
 - String is just page numbers, not full addresses
 - Repeated access to the same page does not cause a page fault

PAGE AND FRAME REPLACEMENT ALGORITHMS

In all our examples, the reference string is

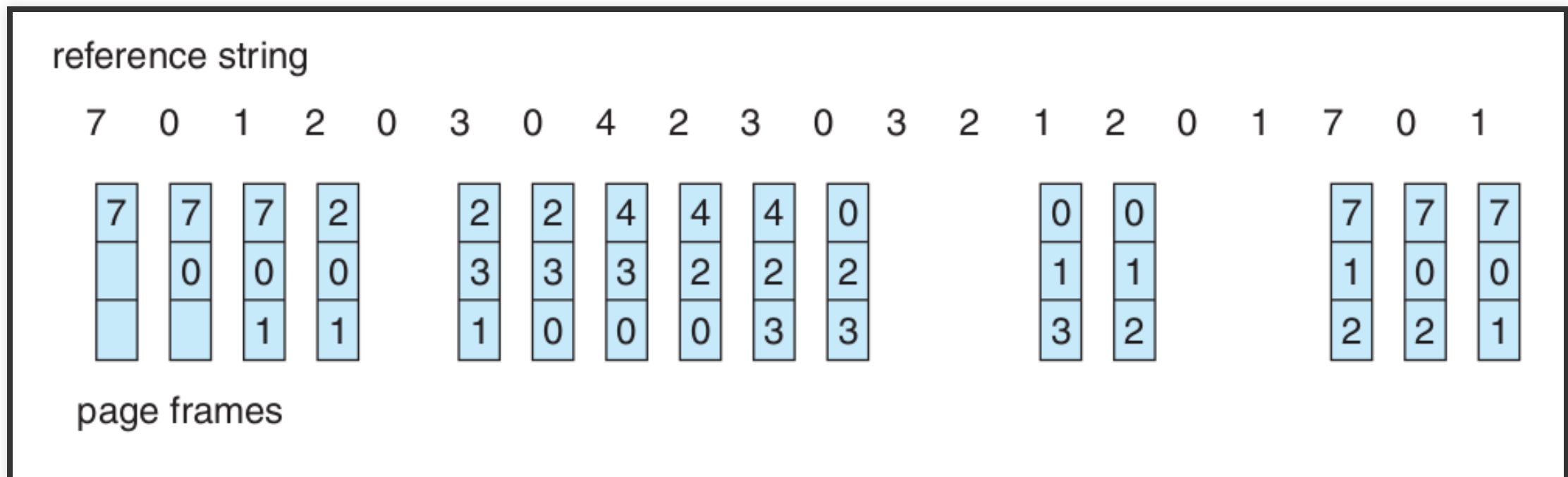
7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 0, 3, 2, 1, 2, 0, 1, 7, 0, 1

GRAPH OF PAGE FAULTS VS THE NUMBER OF FRAMES



FIRST-IN-FIRST-OUT (FIFO) ALGORITHM

- 3 frames (3 pages can be in memory at a time per process)

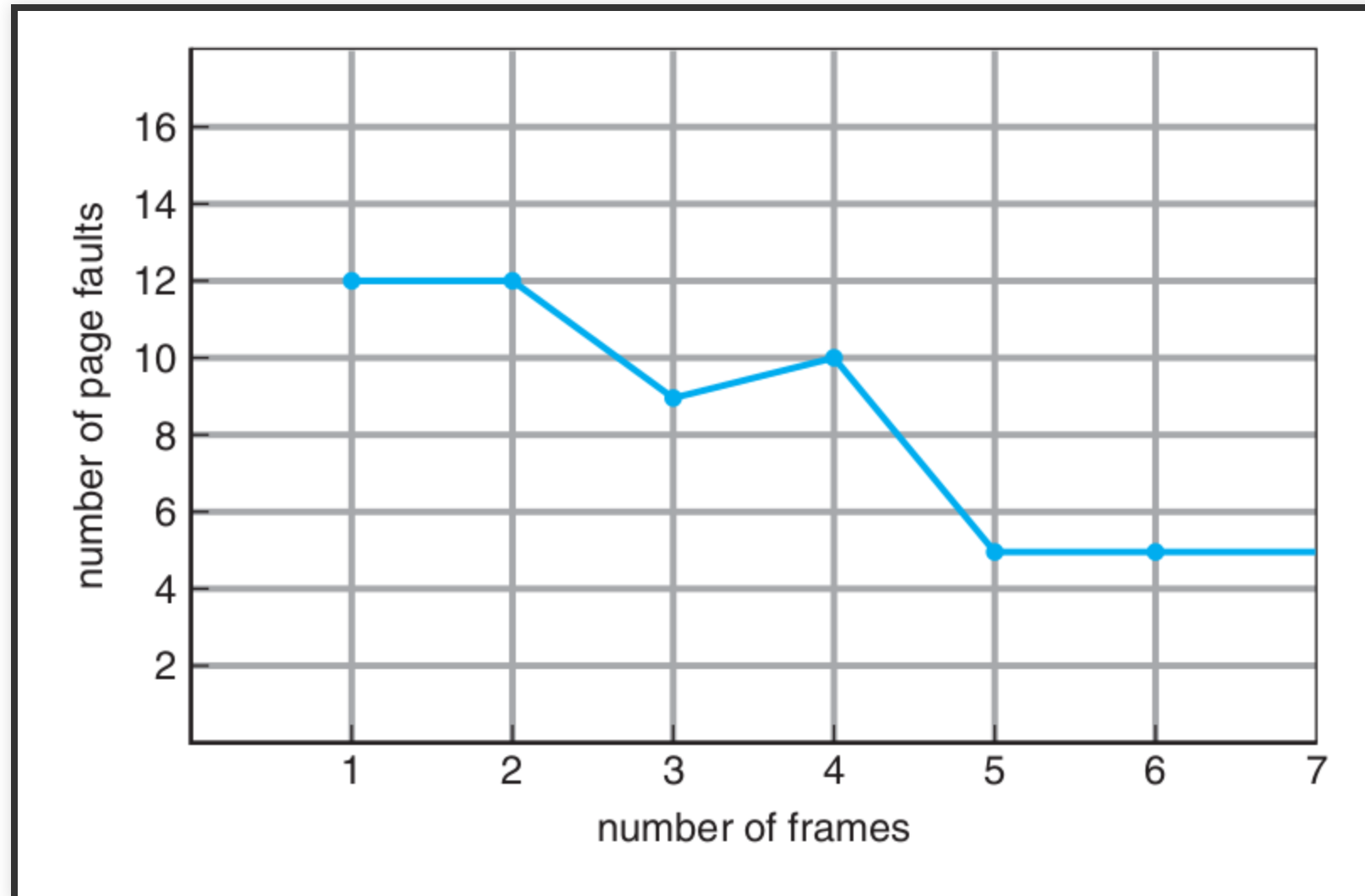


- 15 page faults

FIFO ALGORITHM

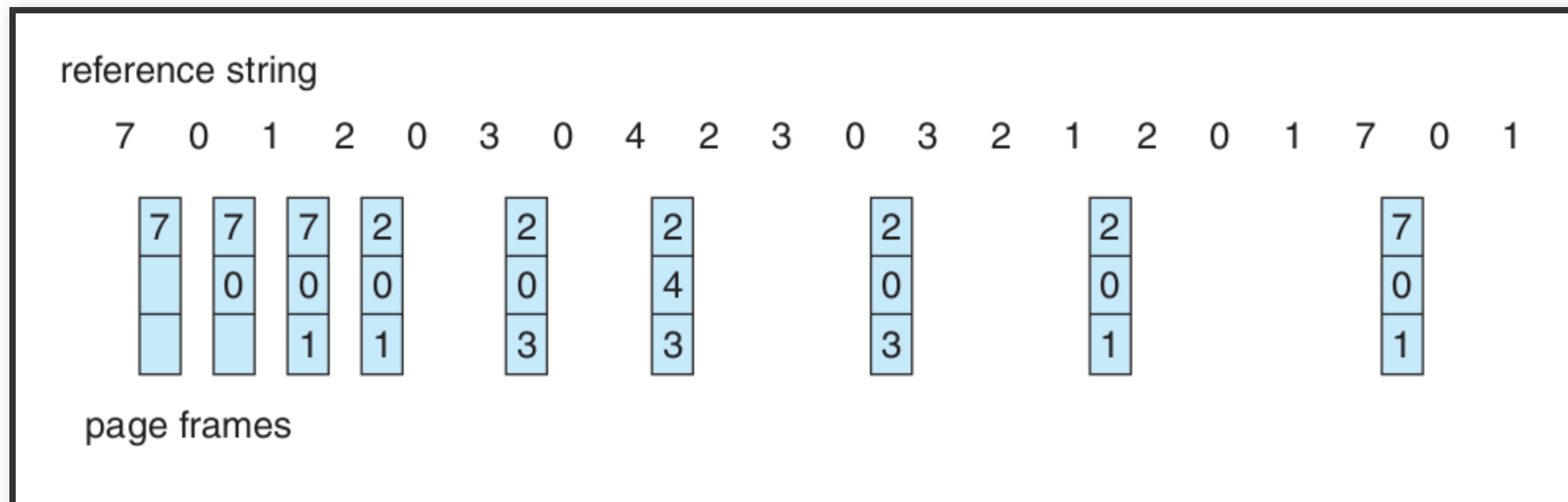
- Can vary by reference string: consider
1,2,3,4,1,2,5,1,2,3,4,5
 - Adding more frames can cause more page faults!
 - Belady's Anomaly
- How to track ages of pages?
 - Just use a FIFO queue

FIFO - BELADY'S ANOMALY



OPTIMAL PAGE REPLACEMENT

- Replace page that will not be used for longest period of time
 - 9 is optimal

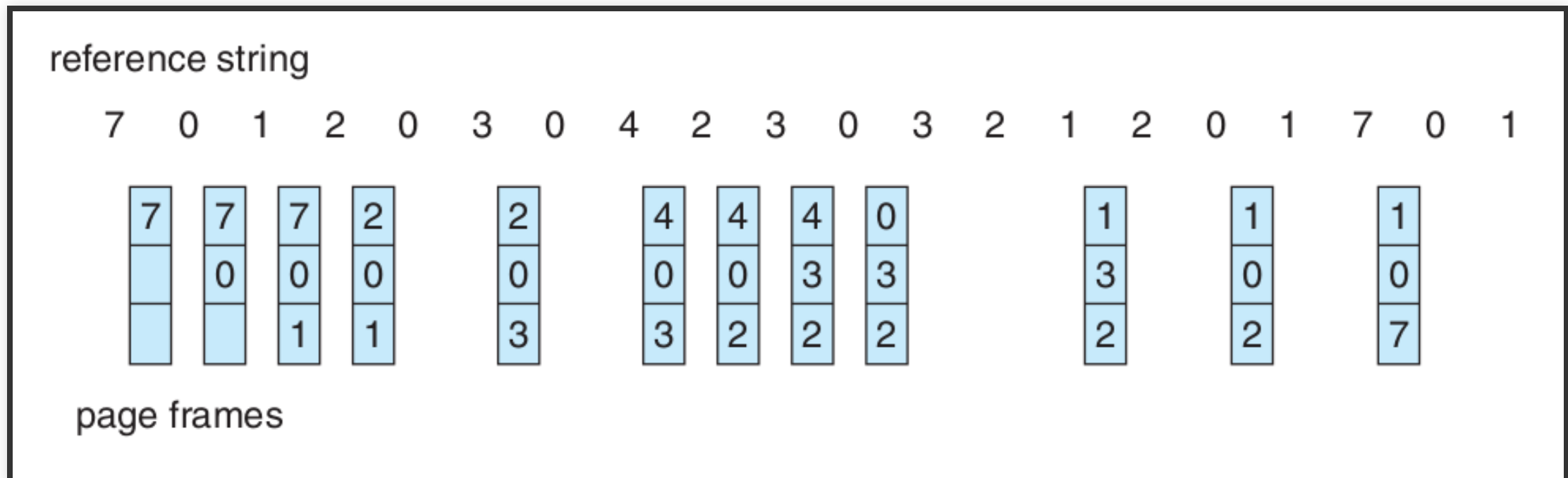


OPTIMAL PAGE REPLACEMENT

- How do you know this?
 - Can't read the future
- Used for measuring how well your algorithm performs

LEAST RECENTLY USED (LRU)

- 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



LEAST RECENTLY USED (LRU)

- 12 faults – better than FIFO but worse than OPT
- Generally good algorithm and frequently used

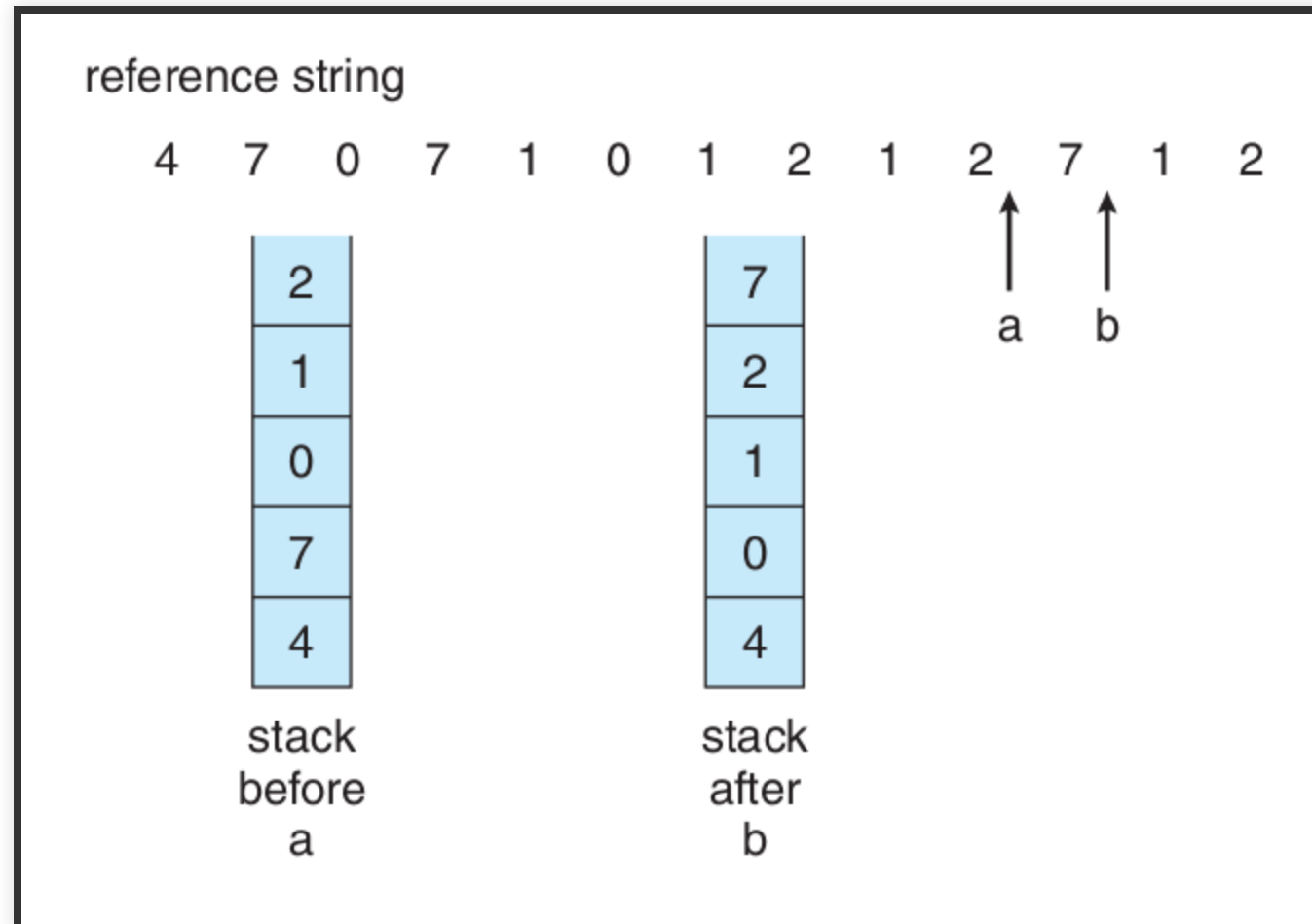
LRU IMPLEMENTATION

- 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

LRU IMPLEMENTATION

- 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 THE MOST RECENT PAGE REFERENCES



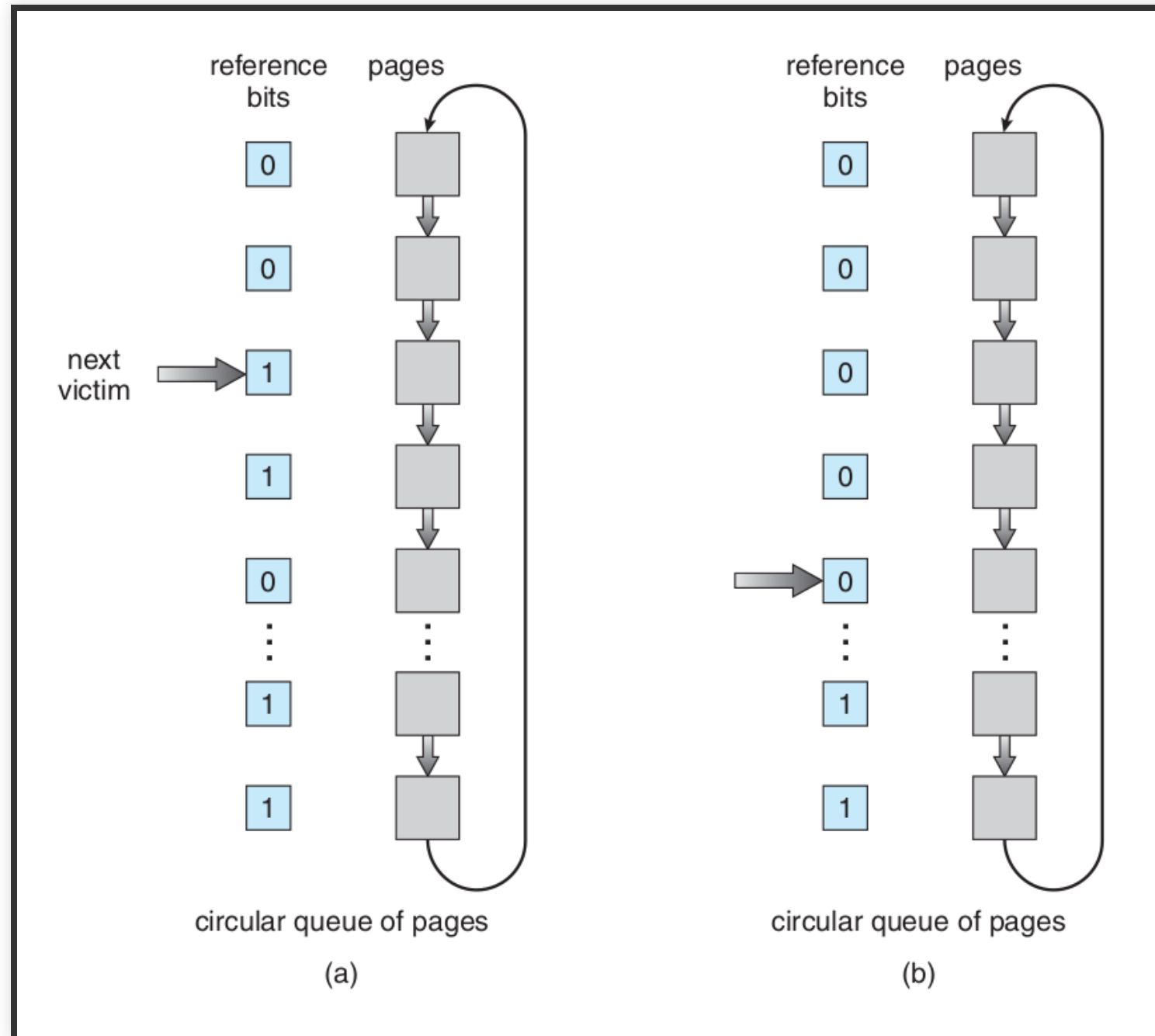
LRU APPROXIMATION ALGORITHMS

- 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

LRU APPROXIMATION ALGORITHMS

- 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

SECOND-CHANCE (CLOCK)



ENHANCED SECOND-CHANCE

Consider reference bit and modify bit as ordered pair

1. (0, 0) neither recently used nor modified—best page to replace
2. (0, 1) not recently used but modified—not quite as good, because the page will need to be written out before replacement
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 the page will need to be written out to secondary storage before it can be replaced

ENHANCED SECOND-CHANCE

- Each page in single class
- Use clock strategy on classes at a time
- replace the first page encountered in the lowest nonempty class
- may have to scan the circular queue several times before we find a page to be replaced.

COUNTING ALGORITHMS

- Keep a counter of the number of references that have been made to each page
 - Not common
- 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

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

PAGE-BUFFERING ALGORITHMS

- 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

APPLICATIONS AND PAGE REPLACEMENT

- All of these algorithms have OS guessing about future page access
- Some applications have better knowledge – i.e. databases
- Memory intensive applications can cause double buffering
 - OS keeps copy of page in memory as I/O buffer
 - Application keeps page in memory for its own work

APPLICATIONS AND PAGE REPLACEMENT

- Operating system can given direct access to the disk, getting out of the way of the applications
 - Raw disk mode
- Bypasses buffering, locking, etc

ALLOCATION OF FRAMES

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

ALLOCATION OF FRAMES

- 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

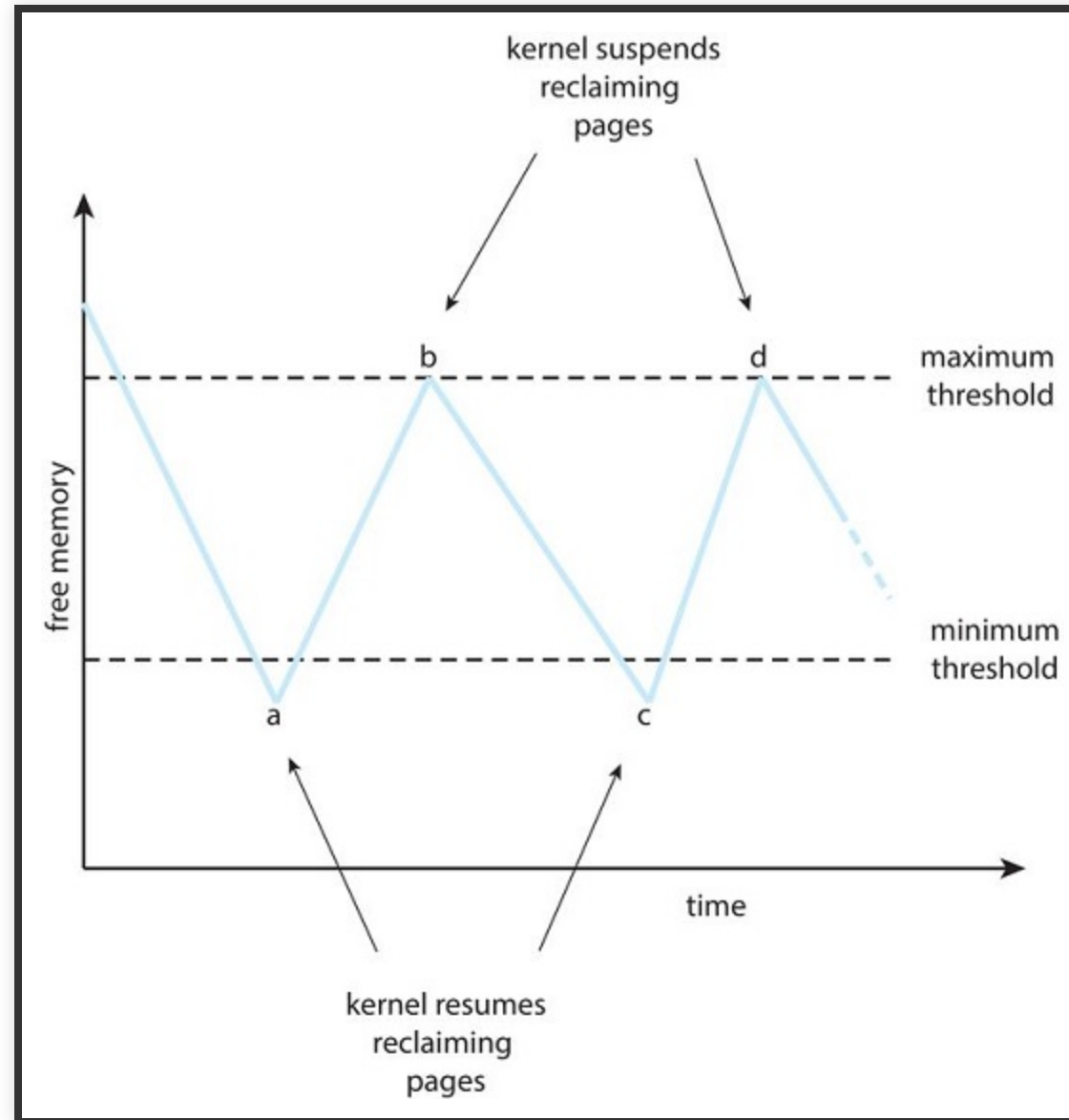
PRIORITY ALLOCATION

- 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

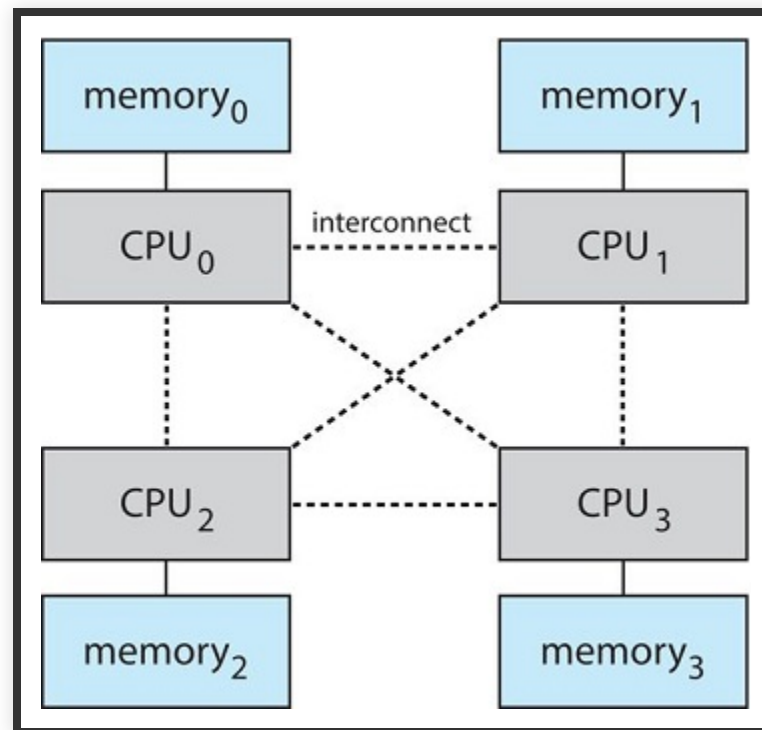
GLOBAL ALLOCATION



NON-UNIFORM MEMORY ACCESS

- So far 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

NON-UNIFORM MEMORY ACCESS



NON-UNIFORM MEMORY ACCESS

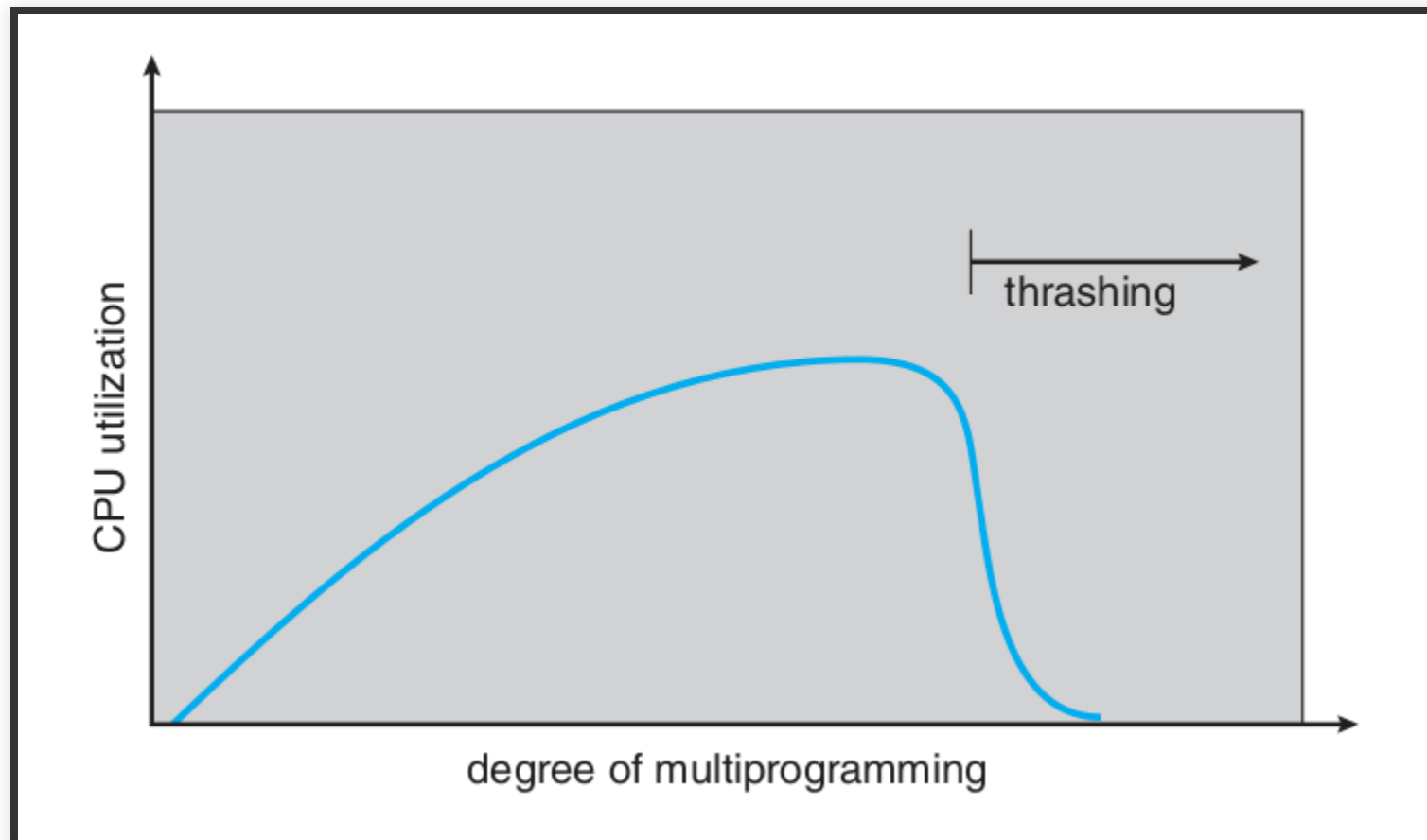
- Optimal performance comes from allocating memory “close to” the CPU on which the thread is scheduled
 - And modifying the scheduler to schedule the thread on the same system board when possible
 - Solved by Solaris by creating Igroups
 - Structure to track CPU / Memory low latency groups
 - Used my schedule and pager
 - When possible schedule all threads of a process and allocate all memory for that process within the Igroup

THRASHING

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
 - ⇒ 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

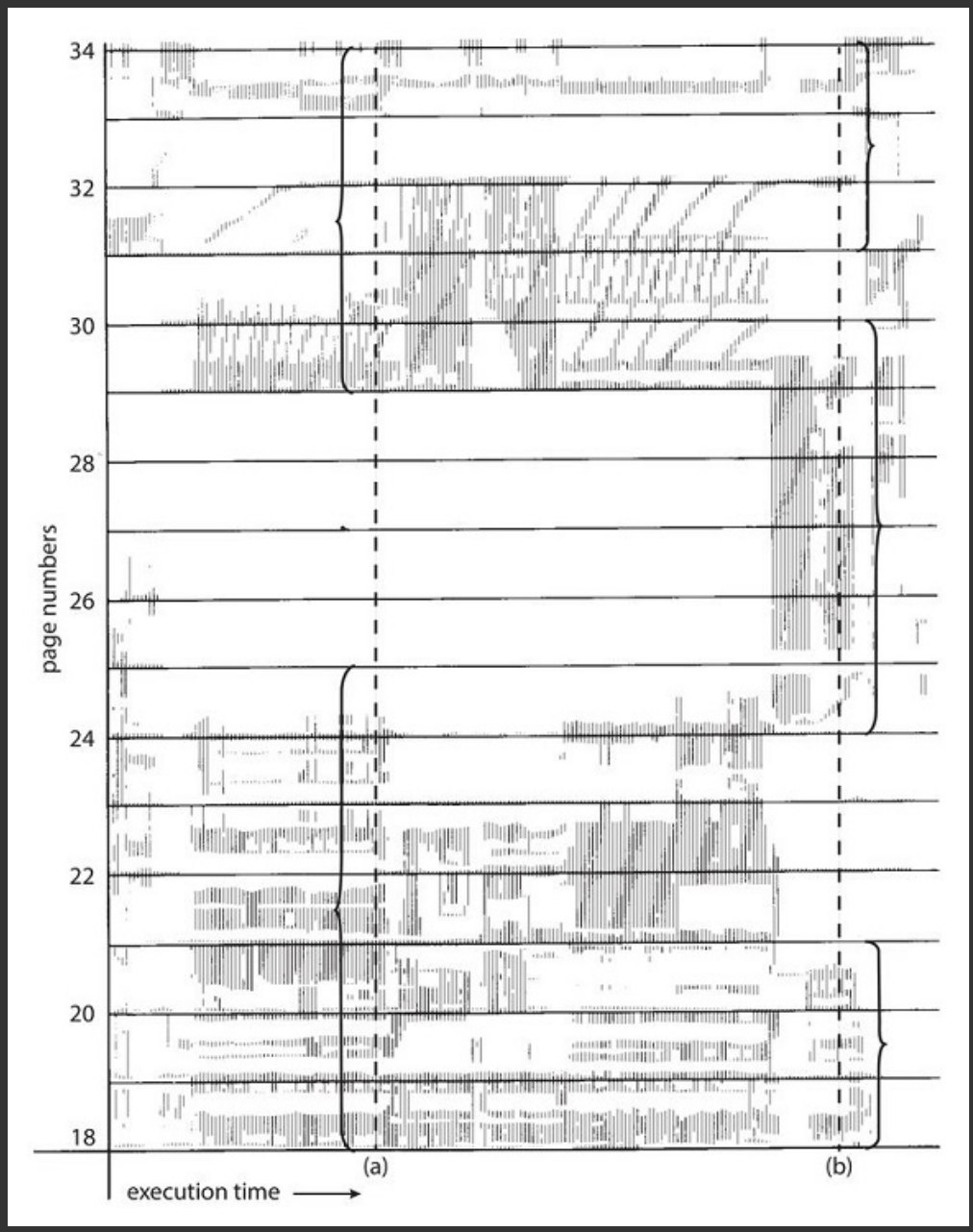
THRASHING



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
 - 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 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

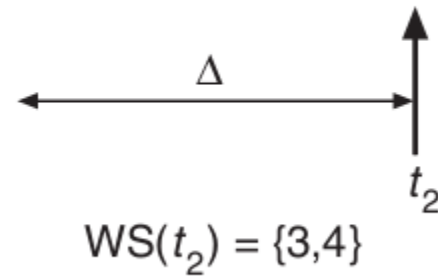
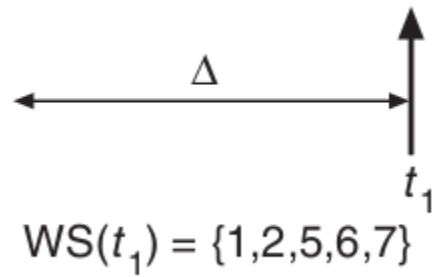
WORKING-SET MODEL

- $D = \sum WSS_i \equiv$ total demand frames
 - Approximation of locality
- if $D > m \Rightarrow$ Thrashing
- Policy if $D > m$, then suspend or swap out one of the processes

WORKING-SET MODEL

page reference table

... 2 6 1 5 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 ...



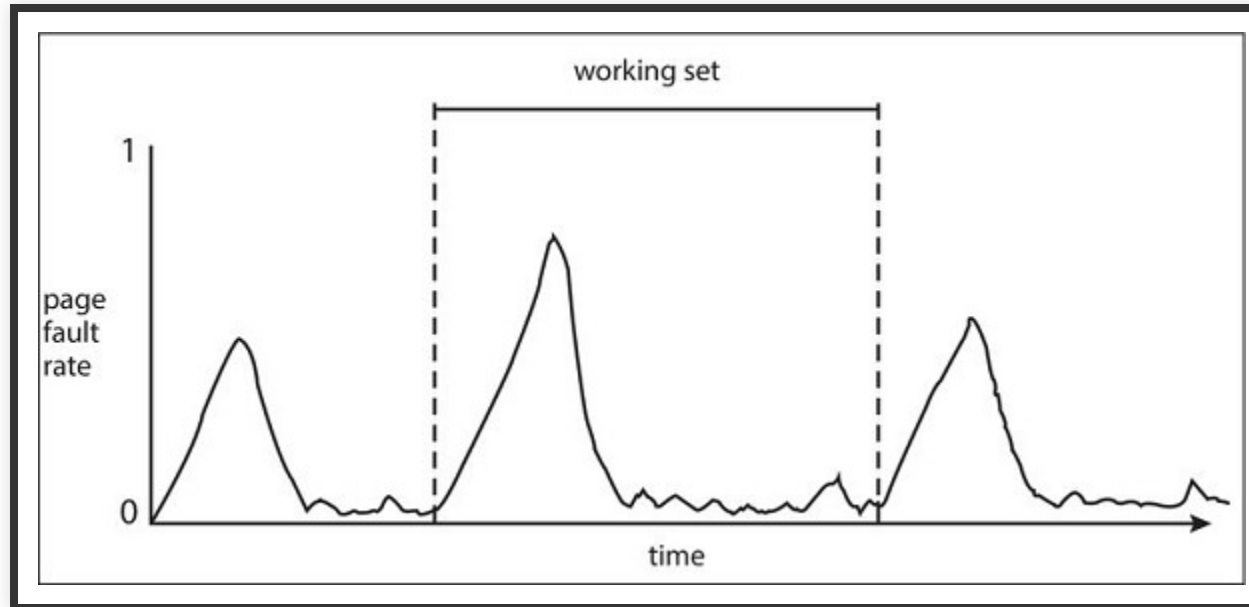
KEEPING TRACK OF WORKING SET

- 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

KEEPING TRACK OF WORKING SET

- Why is this not completely accurate?
- Improvement = 10 bits and interrupt every 1000 time units

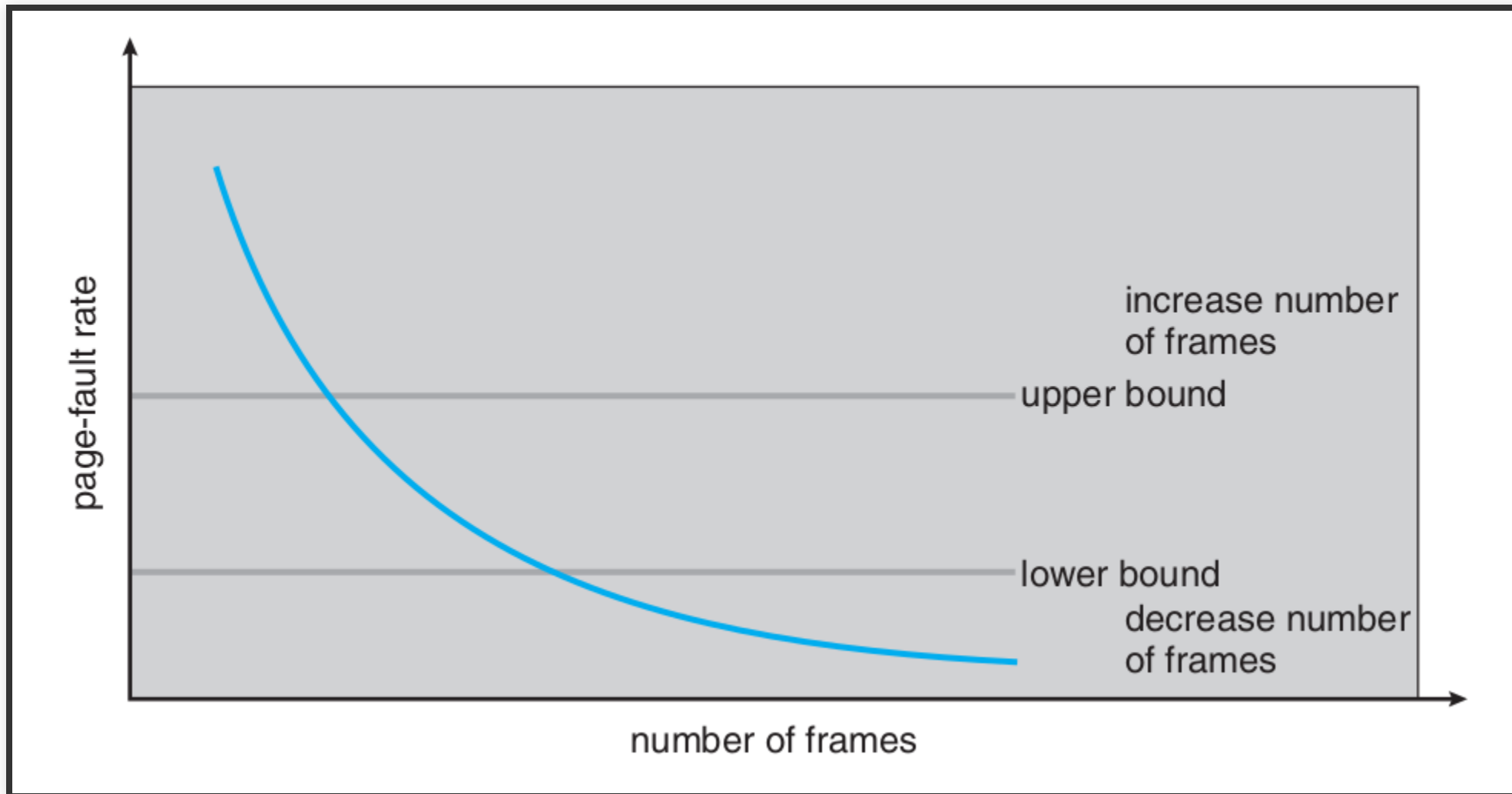
WORKING SET AND PAGE-FAULT RATES



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

PAGE-FAULT FREQUENCY

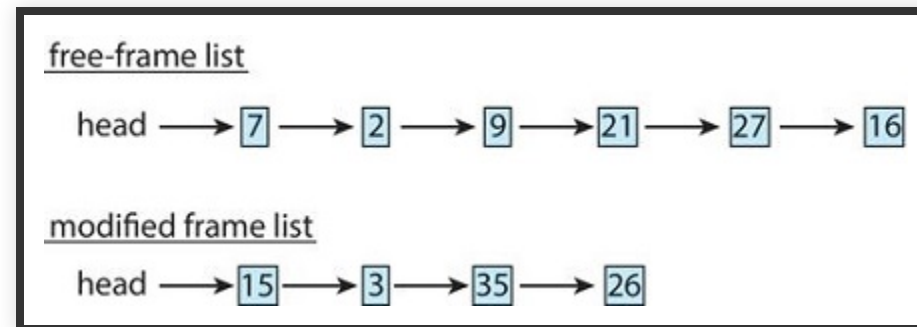


MEMORY COMPRESSION

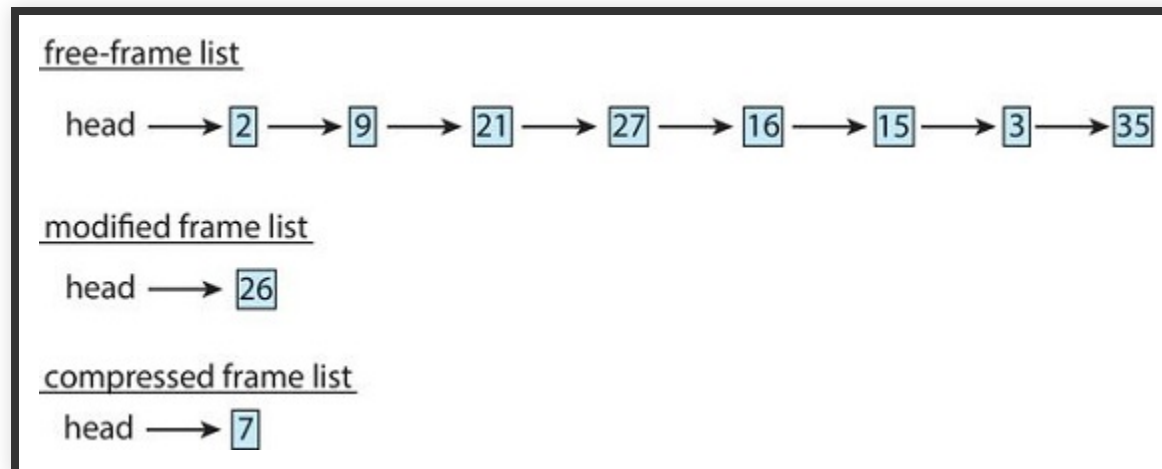
MEMORY COMPRESSION

- An alternative to paging
- rather than paging out modified frames to swap space, we compress several frames into a single frame

MEMORY COMPRESSION



MEMORY COMPRESSION



ALLOCATING KERNEL MEMORY

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
 - I.e. for device I/O

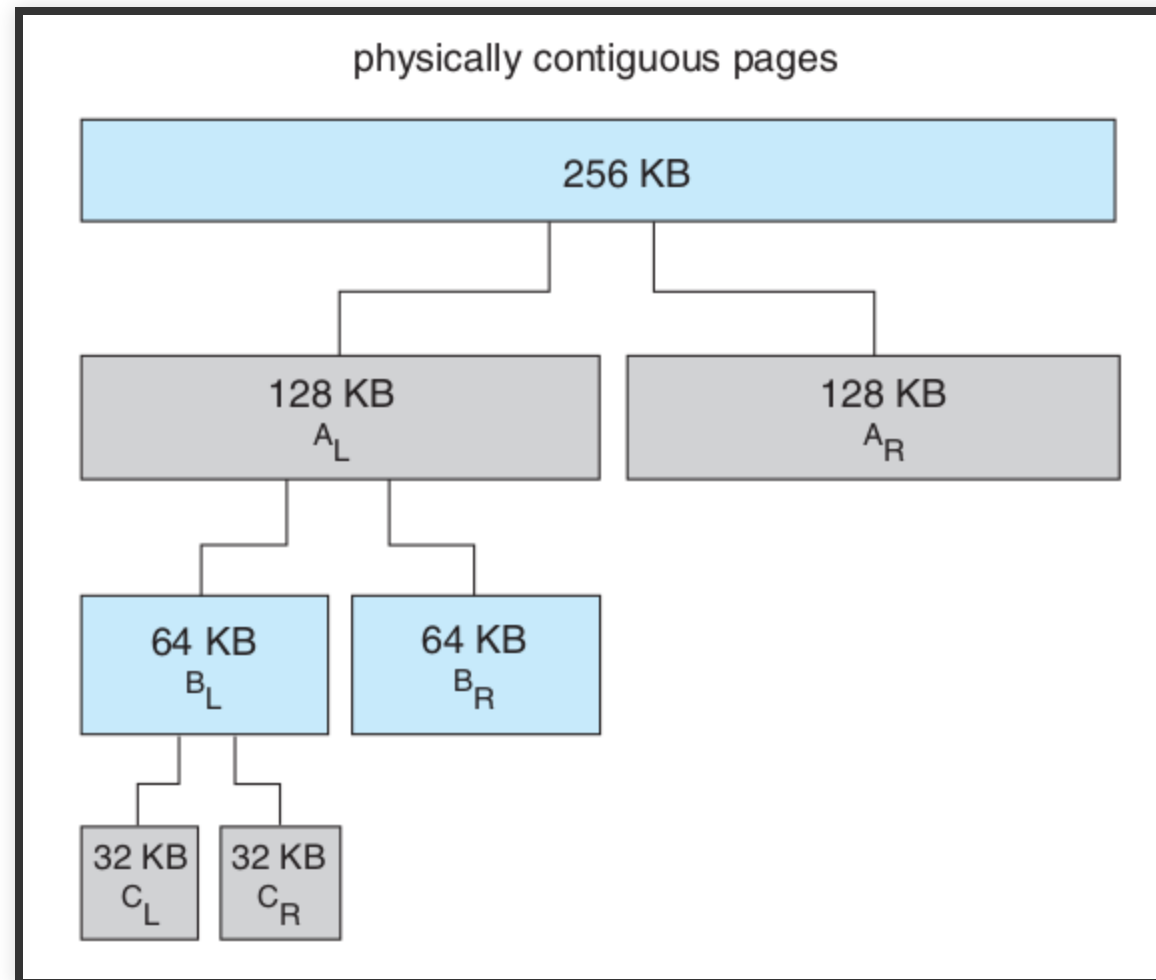
BUDDY SYSTEM

- Allocates memory from fixed-size segment consisting of physically-contiguous pages
- Memory allocated using power-of-2 allocator
 - Satisfies requests in units sized as power of 2
 - Request rounded up to next highest power of 2
 - When smaller allocation needed than is available, current chunk split into two buddies of next-lower power of 2
 - Continue until appropriate sized chunk available

BUDDY SYSTEM

- For example, assume 256KB chunk available, kernel requests 21KB
 - Split into A_L and A_R of 128KB each
 - One further divided into B_L and B_R of 64KB
 - One further into C_L and C_R of 32KB each – one used to satisfy request
- Advantage – quickly coalesce unused chunks into larger chunk
- Disadvantage - fragmentation

BUDDY SYSTEM ALLOCATOR



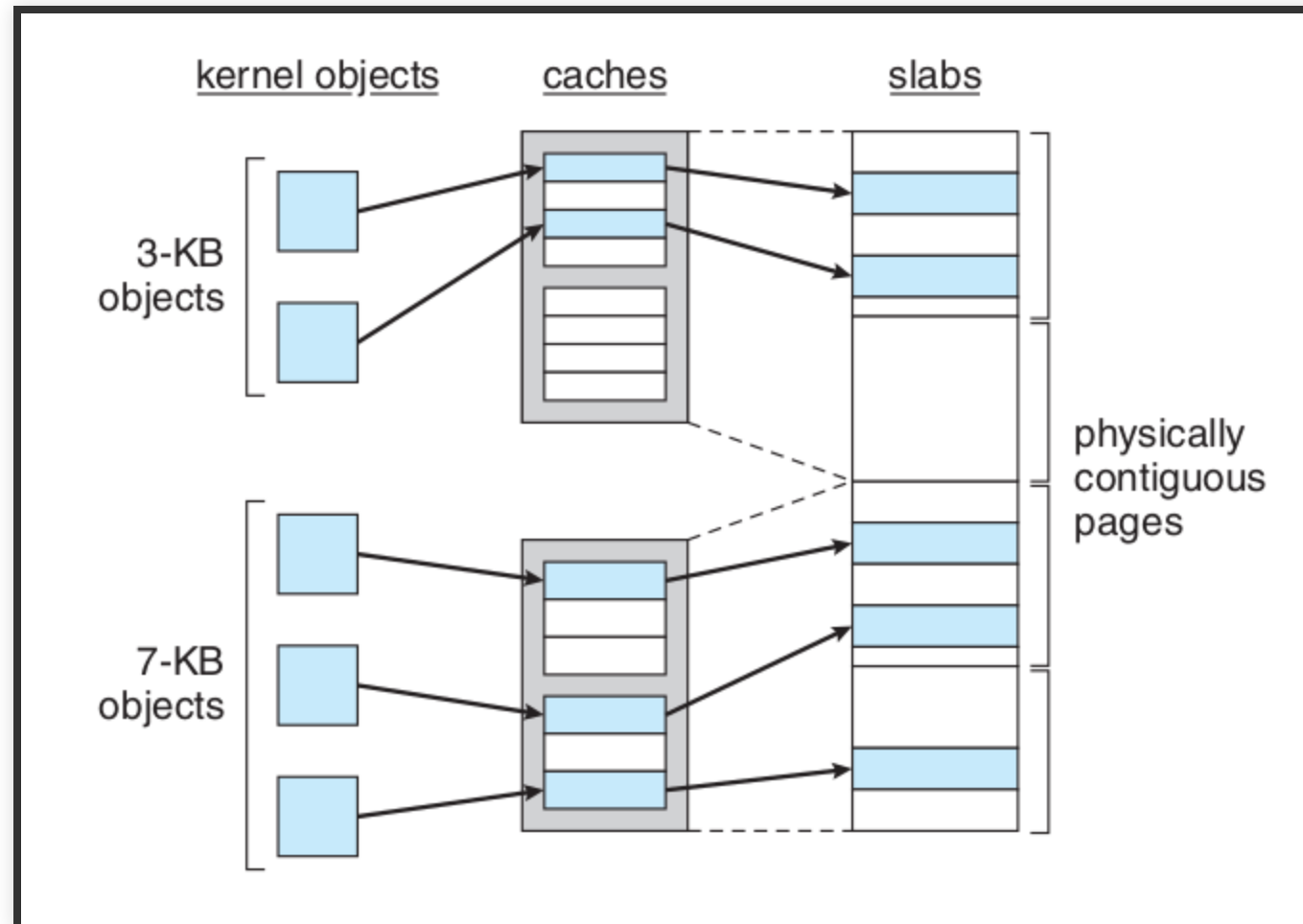
SLAB ALLOCATOR

- Alternate strategy
- Slab is one or more physically contiguous pages
- Cache consists of one or more slabs
- Single cache for each unique kernel data structure
 - Each cache filled with objects – instantiations of the data structure

SLAB ALLOCATOR

- When cache created, filled with objects marked as free
- When structures stored, objects marked as used
- If slab is full of used objects, next object allocated from empty slab
 - If no empty slabs, new slab allocated
- Benefits include no fragmentation, fast memory request satisfaction

SLAB ALLOCATION



SLAB ALLOCATION

In Linux, a slab may be in one of three possible states:

1. Full. All objects in the slab are marked as used.
2. Empty. All objects in the slab are marked as free.
3. Partial. The slab consists of both used and free objects.

OTHER CONSIDERATIONS

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

PREPAGING

- 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

PAGE SIZE

- Sometimes OS designers have a choice
 - Especially if running on custom-built CPU
- Always power of 2, usually in the range 2^{12} (4,096 bytes) to 2^{22} (4,194,304 bytes)
- On average, growing over time

PAGE SIZE

- Page size selection must take into consideration:
 - Fragmentation
 - Page table size
 - Resolution
 - I/O overhead
 - Number of page faults
 - Locality
 - TLB size and effectiveness

TLB REACH

- TLB Reach - The amount of memory accessible from the TLB
- $\text{TLB Reach} = (\text{TLB Size}) * (\text{Page Size})$
- Ideally, the working set of each process is stored in the TLB
 - Otherwise there is a high degree of page faults

TLB REACH

- 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

PROGRAM STRUCTURE

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

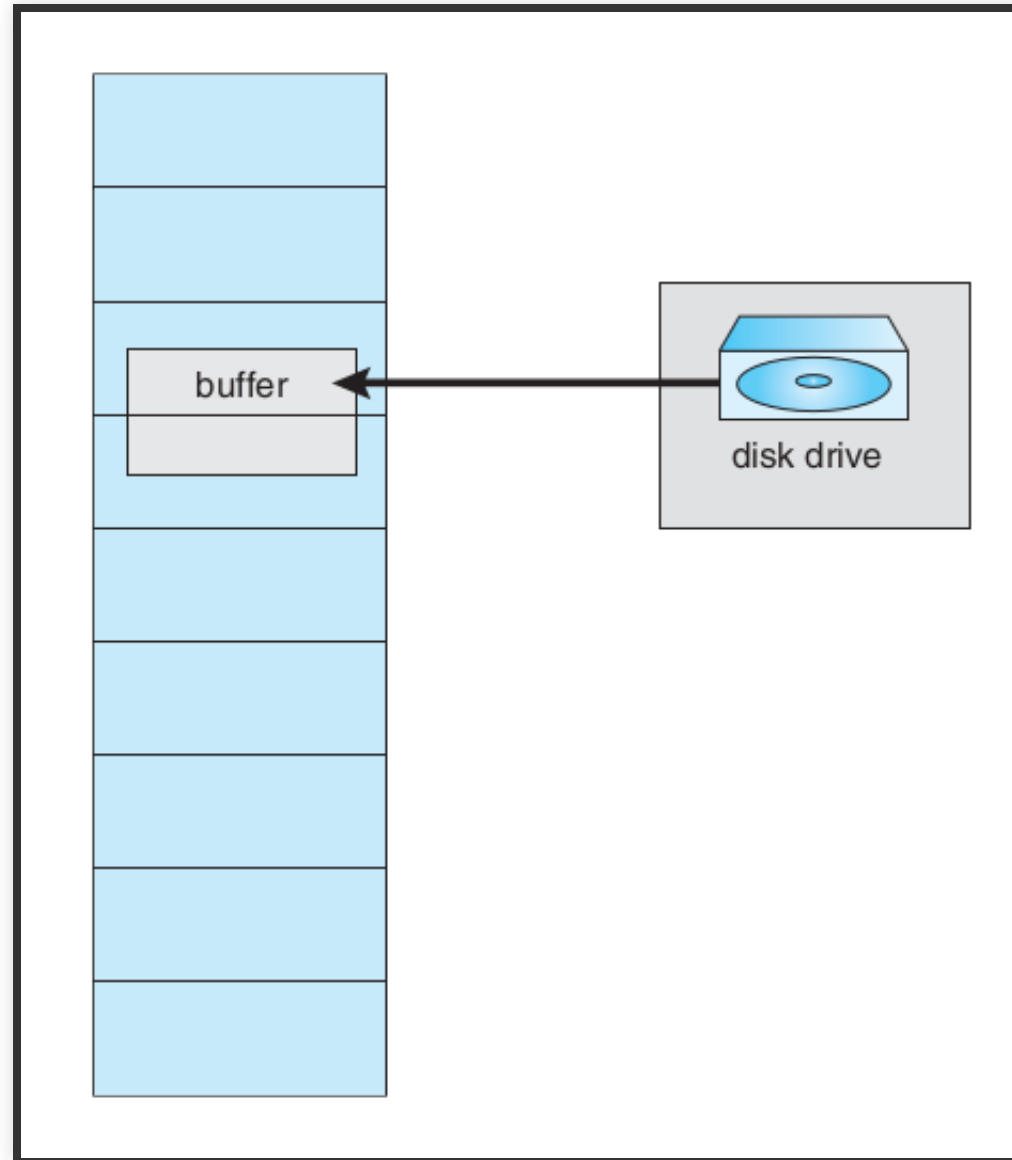
```
int i, j;  
int[128][128] data;  
  
for (j = 0; j < 128; j++)  
    for (i = 0; i < 128; i++)  
        data[i][j] = 0;
```

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

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

IO AND MEMORY

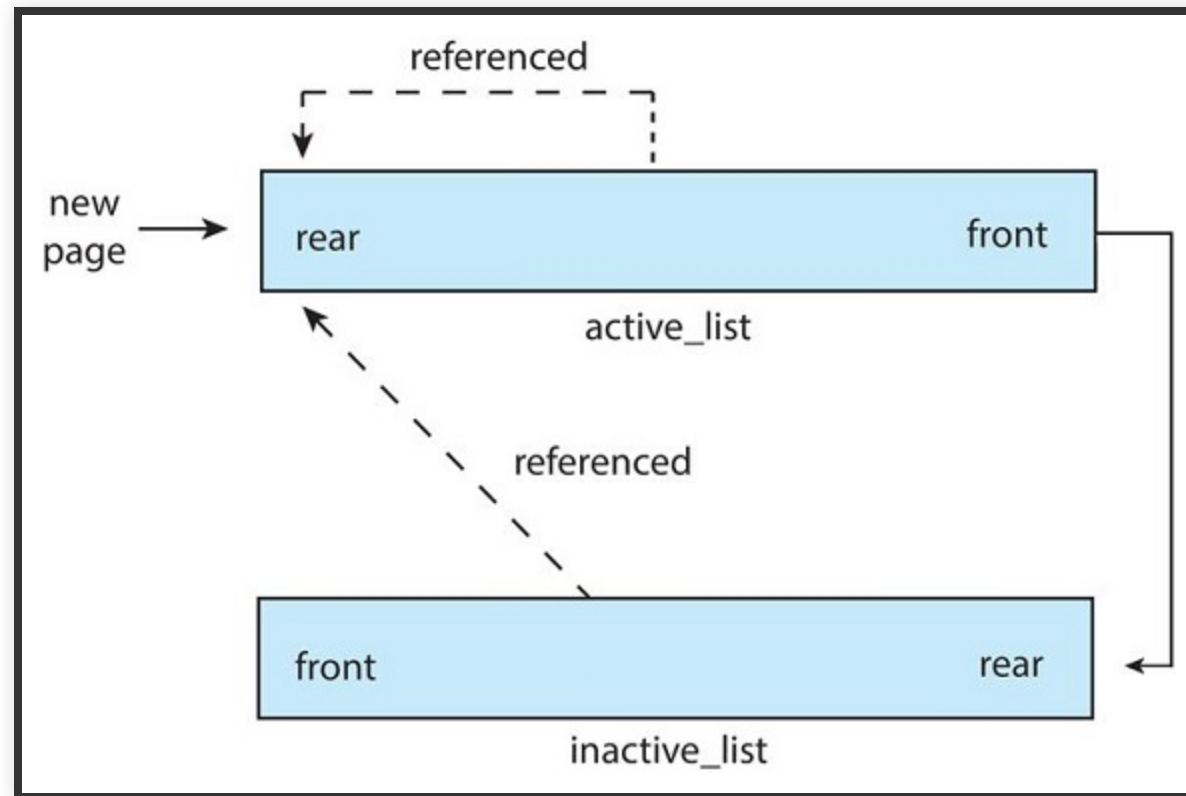


OPERATING-SYSTEM EXAMPLES

LINUX

- Uses demand paging, allocating pages from a list of free frames.
- Global page-replacement policy similar to the LRU-approximation clock algorithm
- Two types of page lists: an `active_list` and an `inactive_list`.
 - The `active_list` contains the pages that are considered in use
 - The `inactive_list` contains pages that have not recently been referenced and are eligible to be reclaimed.

LINUX



WINDOWS

- 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

WINDOWS

- 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

QUESTIONS

BONUS



Exam question number 7: **Virtual Memory**