# 10. Virtual Memory

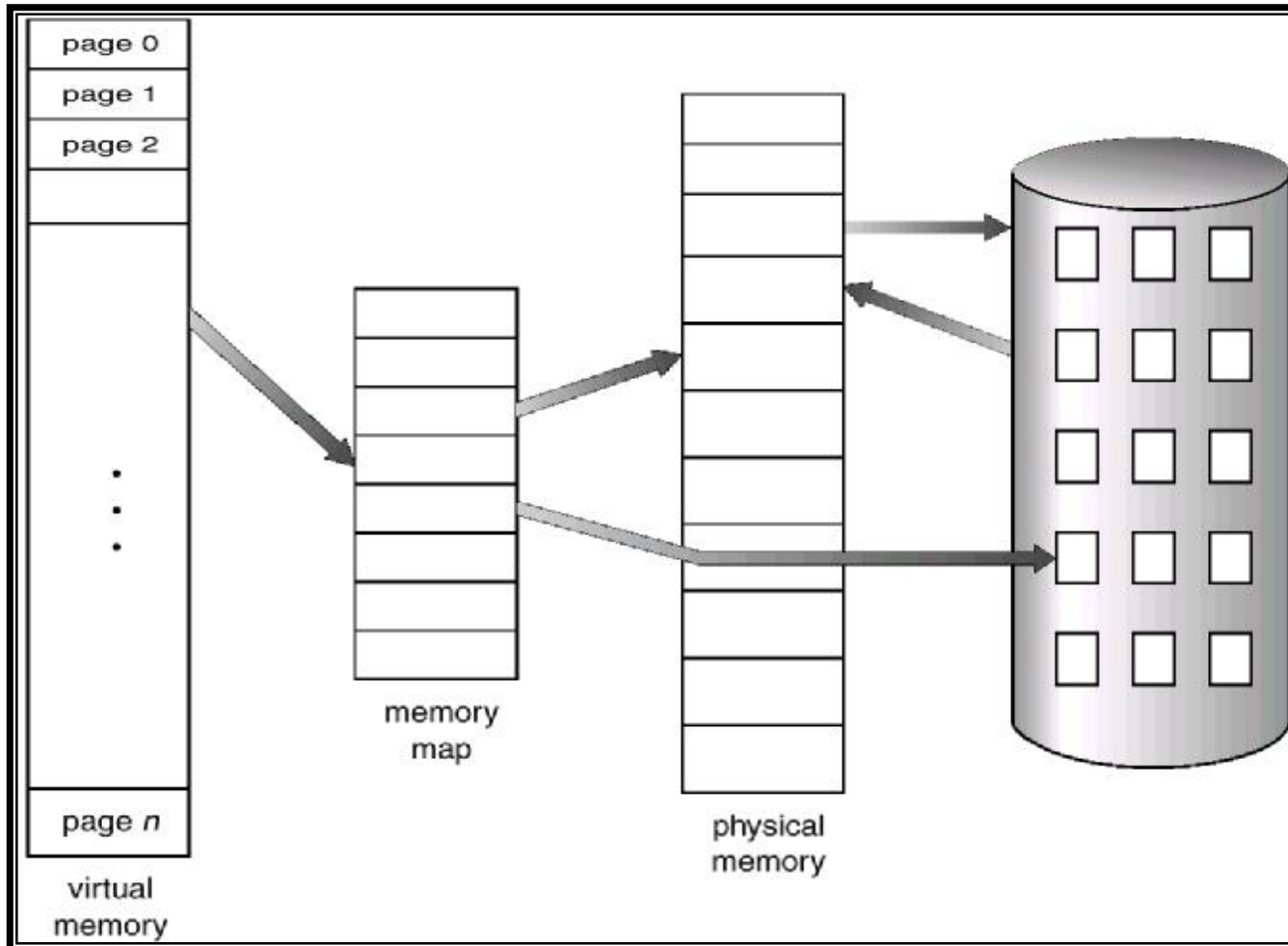*Sungyoung Lee*

*College of Engineering*

*KyungHee University*

# Contents

- *Background*
- *Demand Paging*
- *Process Creation*
- *Page Replacement*
- *Allocation of Frames*
- *Thrashing*
- *Operating System Examples*

# Background

- **n** **Virtual memory** – separation of user logical memory from physical memory
    - ü Only part of the program needs to be in memory for execution
    - ü Logical address space can therefore be much larger than physical address space
    - ü Allows address spaces to be shared by several processes
    - ü Allows for more efficient process creation

- **n** Virtual memory can be implemented via:
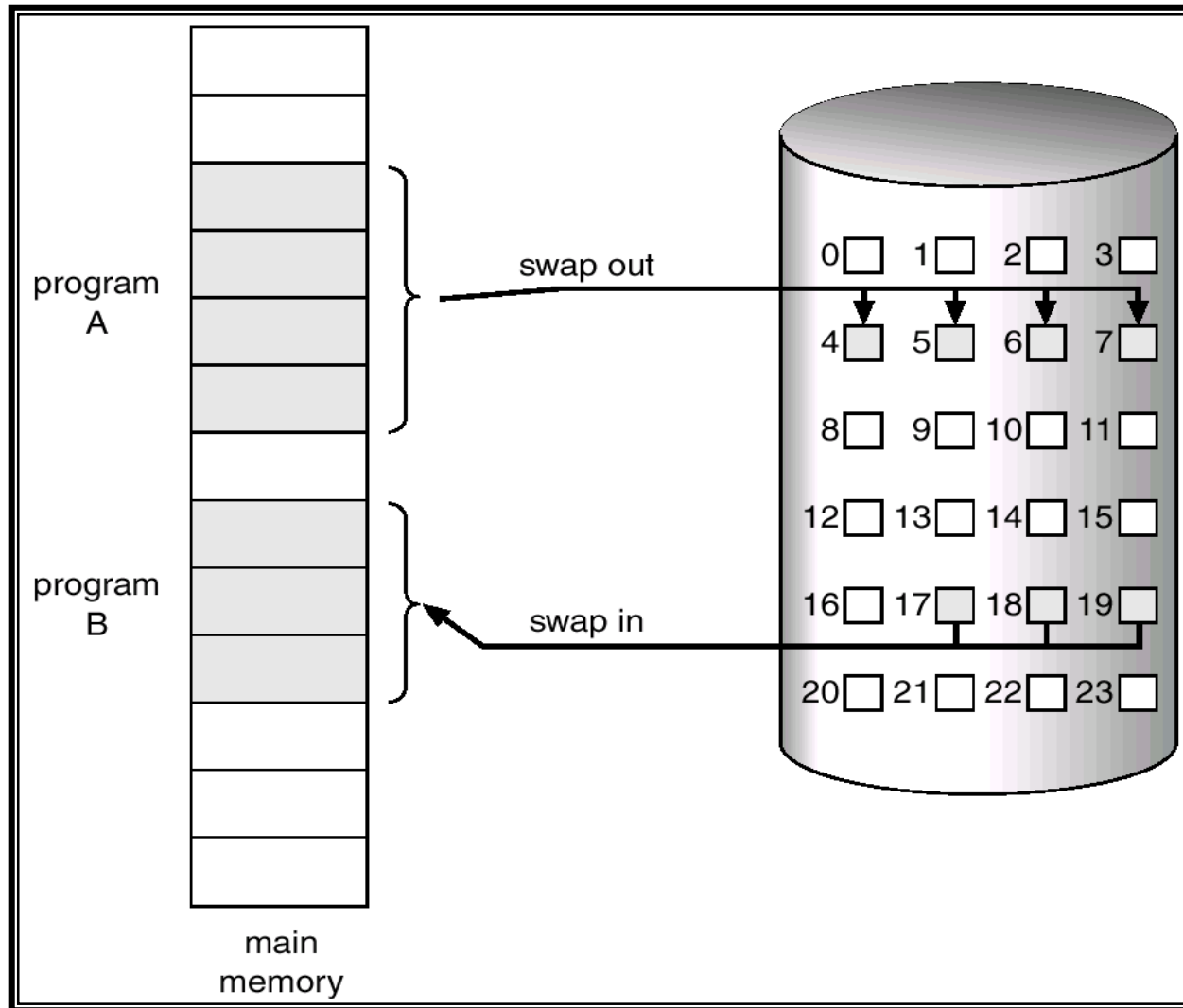    - ü Demand paging
    - ü Demand segmentation

# Demand Paging

**n** Bring a page into memory only when it is needed
- **ü** Less I/O needed
- **ü** Less memory needed
- **ü** Faster response
- **ü** More users

**n** Page is needed $\Rightarrow$ reference to it
- **ü** invalid reference $\Rightarrow$ abort
- **ü** not-in-memory $\Rightarrow$ bring to memory

# Transfer of a Paged Memory to Contiguous Disk Space

- **n** A paging system with (page-level) swapping
- **n** Bring a page into memory only when it is needed
    - ü Cf) swapping: entire process is moved
- **n** OS uses main memory as a (page) cache of all of the data allocated by processes in the system
    - ü Initially, pages are allocated from physical memory frames
    - ü When physical memory fills up, allocating a page requires some other page to be evicted from its physical memory frame
- **n** Evicted pages go to disk (only need to write if they are dirty)
    - ü To a swap file
    - ü Movement of pages between memory/disks is done by the OS
    - ü Transparent to the application

**n** Why does this work? **à** Locality

- **ü** Temporal locality: locations referenced recently tend to be referenced again soon
- **ü** Spatial locality: locations near recently referenced locations are likely to be referenced soon

**n** Locality means paging can be infrequent

- **ü** Once you've paged something in, it will be used many times
- **ü** On average, you use things that are paged in
- **ü** But this depends on many things:
  - **§** Degree of locality in application
  - **§** Page replacement policy
  - **§** Amount of physical memory
  - **§** Application's reference pattern and memory footprint

**n** **Why is this "demand" paging?**

- **ü** When a process first starts up, it has a brand new page table, with all PTE valid bits "false"
    - § No pages are yet mapped to physical memory
- **ü** When the process starts executing:
    - § Instructions immediately fault on both code and data pages
    - § Faults stop when all necessary code/data pages are in memory
    - § Only the code/data that is needed (demanded!!) by process needs to be loaded
    - § What is needed changes over time, of course…
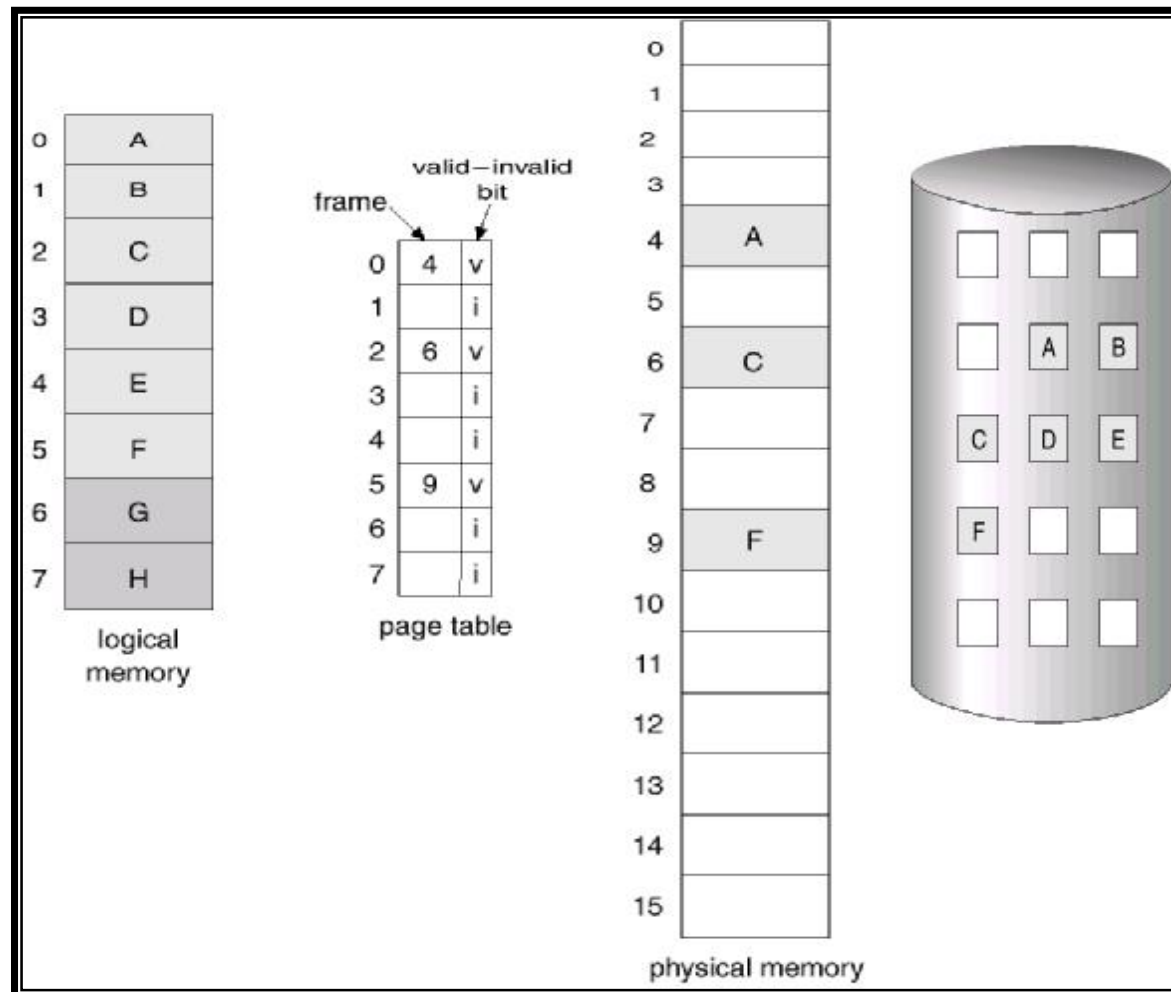
# Valid-Invalid Bit

n  With each page table entry a valid–invalid bit is associated
   (1 $\Rightarrow$ in-memory, 0 $\Rightarrow$ not-in-memory)

n  Initially valid–invalid but is set to 0 on all entries

n  Example of a page table snapshot

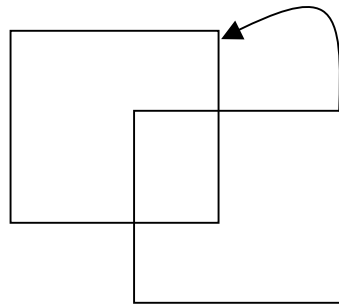| Frame # | valid-invalid bit |
|---------|-------------------|
|         | 1                 |
|         | 1                 |
|         | 1                 |
|         | 1                 |
|         | 0                 |
| M       |                   |
|         | 0                 |
|         | 0                 |

page table

n  During address translation, if valid–invalid bit in page table entry is 0
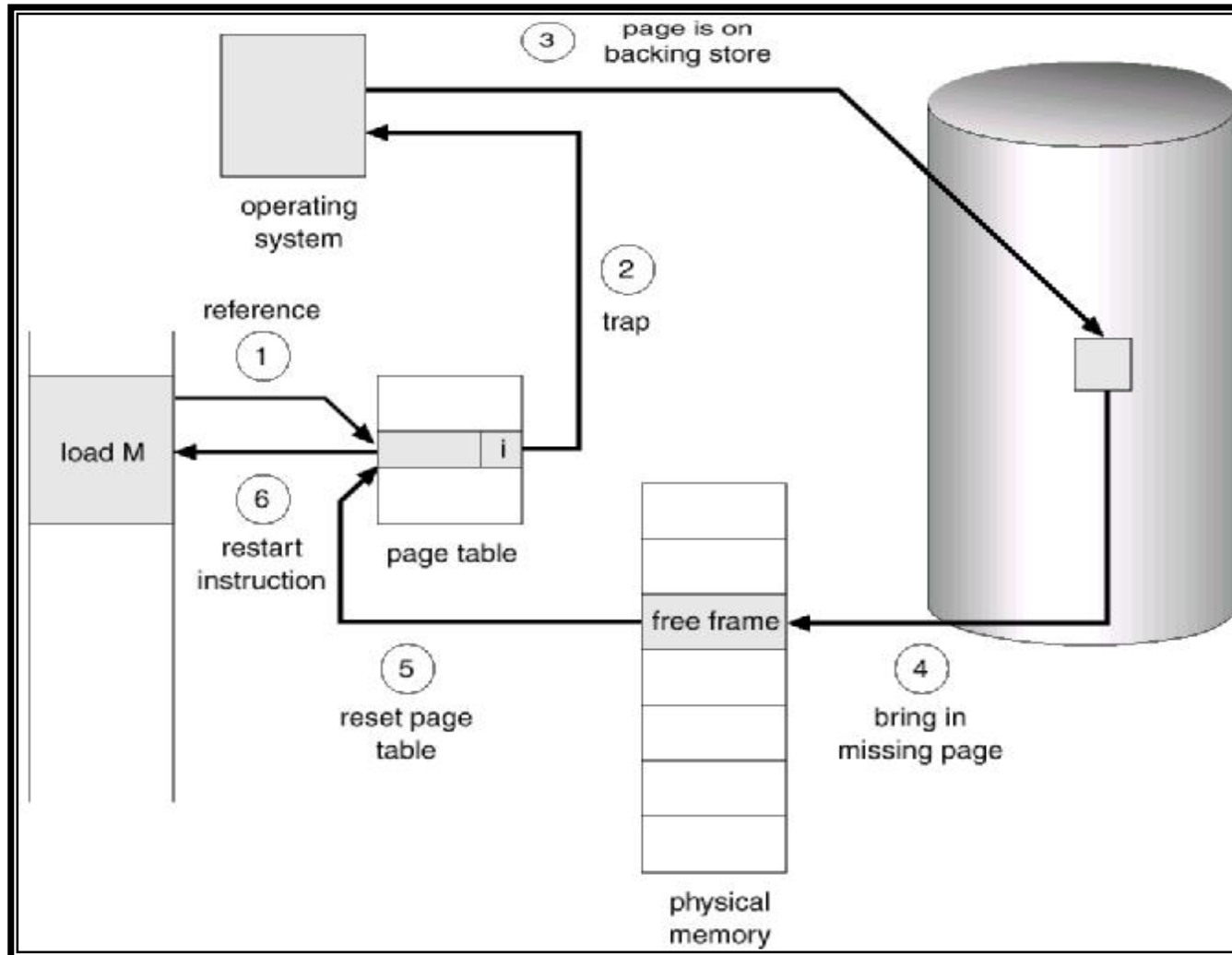   $\Rightarrow$ page fault

# Page Fault

**n** If there is ever a reference to a page, first reference will trap to OS
$\Rightarrow$ page fault

**n** OS looks at another table to decide:

   **ü** Invalid reference $\Rightarrow$ abort

   **ü** Just not in memory

**n** Get empty frame

**n** Swap page into frame

**n** Reset tables, validation bit = 1

**n** Restart instruction (if cannot be restarted?)

   **ü** block move

   **ü** auto increment/decrement location

n  What happens to a process that references a virtual address in a page that has been evicted?

  ü When the page was evicted, the OS sets the PTE as invalid and stores (in PTE) the location of the page in the swap file

  ü When a process accesses the page, the invalid PTE will cause an exception to be thrown

n  The OS will run the page fault handler in response

  ü Handler uses invalid PTE to locate page in swap file

  ü Handler reads page into a physical frame, updates PTE to point to it and to be valid

  ü Handler restarts the faulted process

n  Where does the page that's read in go?

  ü Have to evict something else (page replacement algorithm)

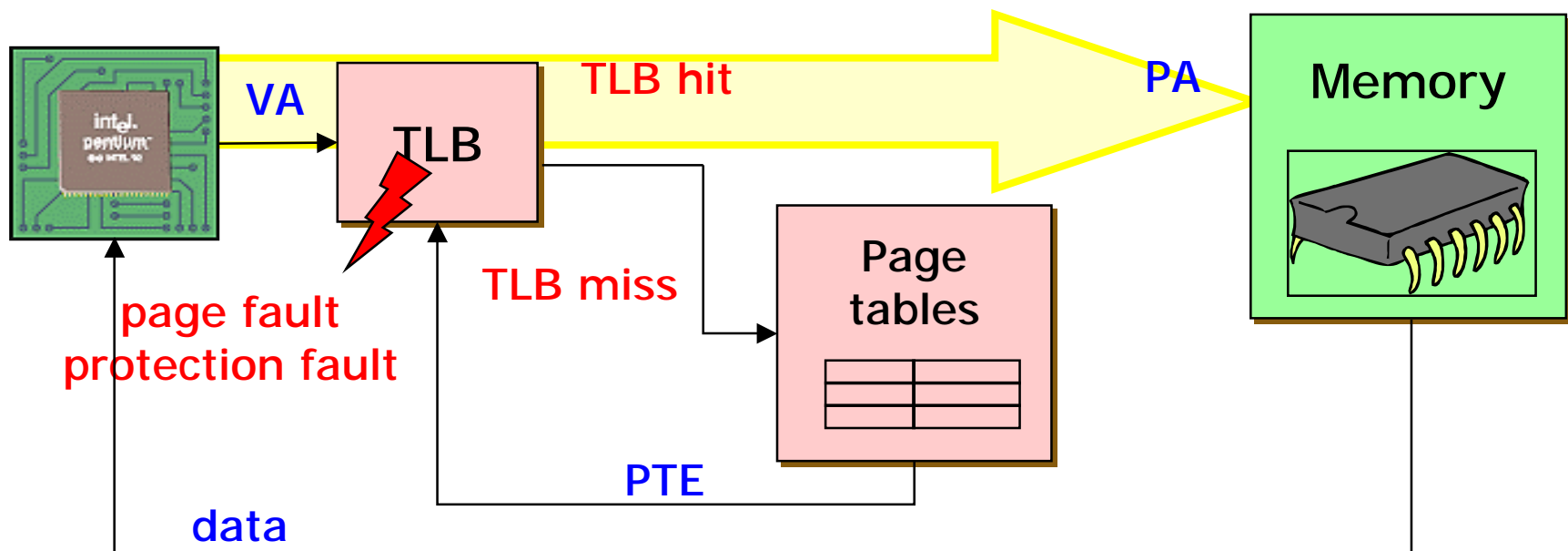  ü OS typically tries to keep a pool of free pages around so that allocations don't inevitably cause evictions

**n** Situation

  **ü** Process is executing on the CPU, and it issues a read to a (virtual) address



TLB hit

VA

TLB

PA

Memory

page fault
protection fault

TLB miss

Page
tables

PTE

data

**n** The common case

- ü The read goes to the TLB in the MMU
- ü TLB does a lookup using the page number of the address
- ü The page number matches, returning a PTE
- ü TLB validates that the PTE protection allows reads
- ü PTE specifies which physical frame holds the page
- ü MMU combines the physical frame and offset into a physical address
- ü MMU then reads from that physical address, returns value to CPU

**n** TLB misses: two possibilities

    ü (1) MMU loads PTE from page table in memory

        § Hardware managed TLB, OS not involved in this step

        § OS has already set up the page tables so that the hardware can access it directly

    ü (2) Trap to the OS

        § Software managed TLB, OS intervenes at this point

        § OS does lookup in page tables, loads PTE into TLB

        § OS returns from exception, TLB continues

    ü At this point, there is a valid PTE for the address in the TLB

**n** TLB misses

    **ü** Page table lookup (by HW or OS) can cause a recursive fault if page table is paged out

        **§** Assuming page tables are in OS virtual address space

        **§** Not a problem if tables are in physical memory

    **ü** When TLB has PTE, it restarts translation

        **§** Common case is that the PTE refers to a valid page in memory

        **§** Uncommon case is that TLB faults again on PTE because of PTE protection bits (e.g., page is invalid)

**n** Page faults

    ü PTE can indicate a protection fault

        § Read/Write/Execute – operation not permitted on page

        § Invalid – virtual page not allocated, or page not in physical memory

    ü TLB traps to the OS (software takes over)

        § Read/Write/Execute – OS usually will send fault back to the process, or might be playing tricks (e.g., copy on write, mapped files)

        § Invalid (Not allocated) – OS sends fault to the process (e.g., segmentation fault)

        § Invalid (Not in physical memory) – OS allocates a frame, reads from disk, and maps PTE to physical frame

# What happens if there is no free frame?

n Page replacement – find some page in memory, but not really in use, swap it out

  ü Algorithm
  ü Performance

    § want an algorithm which will result in minimum number of page faults

n Same page may be brought into memory several times

# Performance of Demand Paging

**n** Page Fault Rate $0 \le p \le 1.0$

    ü if $p = 0$ no page faults

    ü if $p = 1$, every reference is a fault

**n** Effective Access Time (EAT)

$$EAT = (1 - p) \text{ x memory access}$$
$$+ \, p \text{ x (page fault overhead}$$
$$+ \, [\text{swap page out }]$$
$$+ \, \text{swap page in}$$
$$+ \, \text{restart overhead})$$

# Demand Paging Example

n  Memory access time = 1 microsecond

n  50% of the time the page that is being replaced has been modified and therefore needs to be swapped out

n  Swap Page Time = 10 msec = 10,000 usec

$$EAT = (1 - p) \times 1 + p \times (15000)$$
$$= 1 + 14999P \quad \text{(in usec)}$$

# Process Creation

**n** Virtual memory allows other benefits during process creation:

  ü Copy-on-Write

  ü Memory-Mapped Files

# Copy-on-Write

- n Copy-on-Write (COW) allows both parent and child processes to initially *share* the same pages in memory
  - ü If either process modifies a shared page, only then is the page copied

- n COW allows more efficient process creation as only modified pages are copied

- n Free pages are allocated from a *pool* of zeroed-out pages

# Copy-On-Write

**n** Process creation

- ü requires copying the entire address space of the parent process to the child process
- ü Very slow and inefficient!

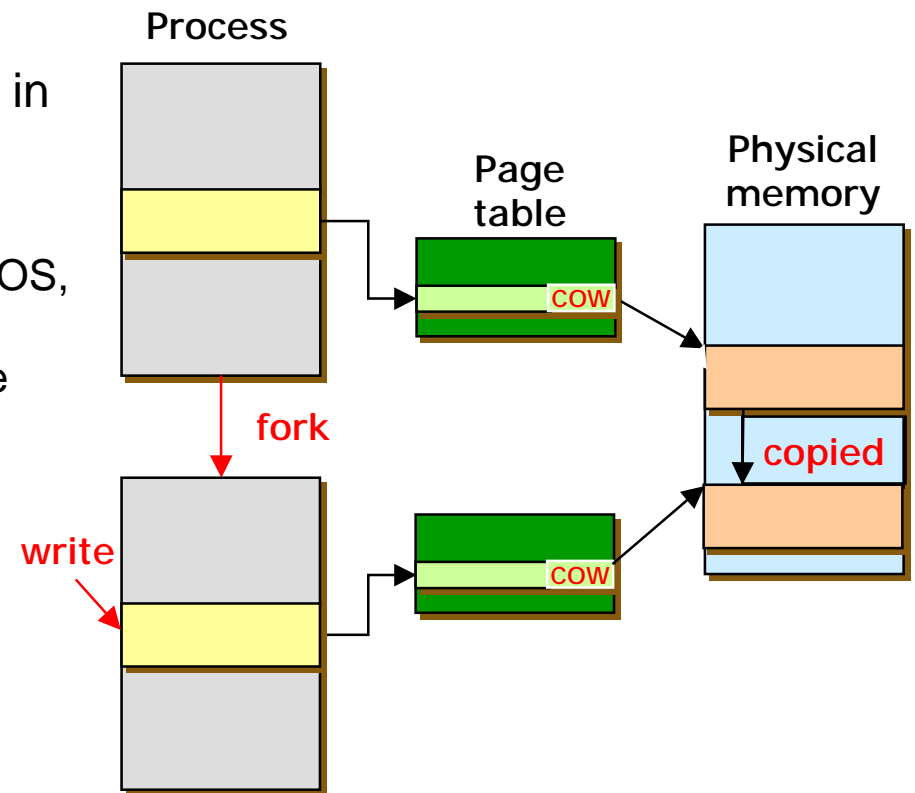**n** Solution 1: Use threads

- ü Sharing address space is free

**n** Solution 2: Use vfork() system call

- ü vfork() creates a process that shares the memory address space of its parent
- ü To prevent the parent from overwriting data needed by the child, the parent's execution is blocked until the child exits or executes a new program
- ü Any change by the child is visible to the parent once it resumes
- ü Useful when the child immediately executes exec()

**n** Solution 3: Copy On Write (COW)

    **ü** Instead of copying all pages, create shared mappings of parent pages in child address space.

    **ü** Shared pages are protected as read-only in child.

        § Reads happen as usual

        § Writes generate a protection fault, trap to OS, and OS copies the page, changes page mapping in client page table, restarts write instruction

**Process**

**Page table**

**Physical memory**
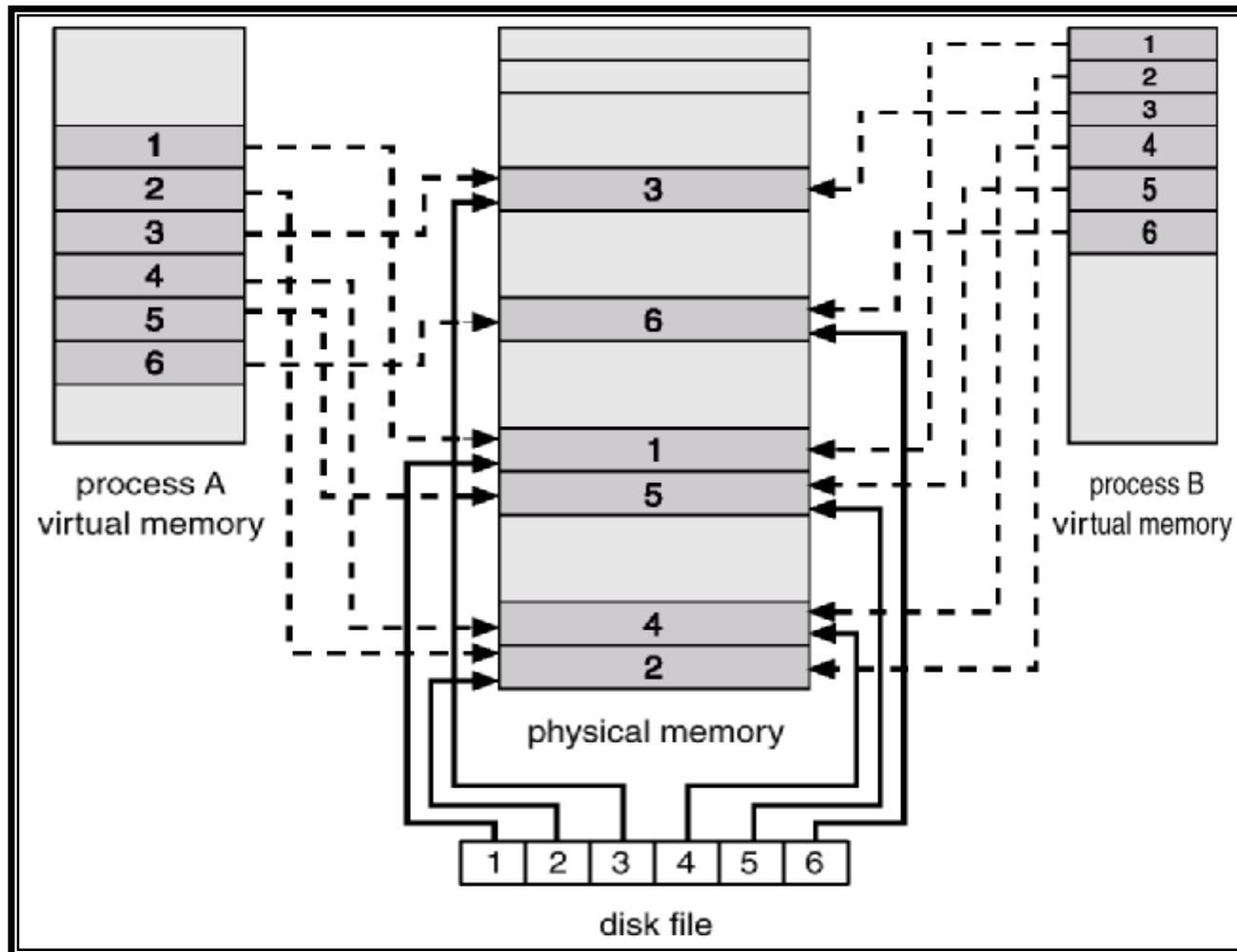
cow

fork

write

cow

copied

# Memory-Mapped Files

n 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

n A file is initially read using demand paging
  ü A page-sized portion of the file is read from the file system into a physical page
  ü Subsequent reads/writes to/from the file are treated as ordinary memory accesses

n Simplifies file access by treating file I/O through memory rather than **read() write()** system calls

n Also allows several processes to map the same file allowing the pages in memory to be shared

# Memory-Mapped Files

n Memory-mapped files
- ü Mapped files enable processes to do file I/O using memory references
    - § Instead of open(), read(), write(), close()
- ü mmap(): bind a file to a virtual memory region
    - § PTEs map virtual addresses to physical frames holding file data
    - § <Virtual address base + N> refers to offset N in file
- ü Initially, all pages in mapped region marked as invalid
    - § OS reads a page from file whenever invalid page is accessed
    - § OS writes a page to file when evicted from physical memory
    - § If page is not dirty, no write needed

**n** Note:
  - ü File is essentially backing store for that region of the virtual address space (instead of using the swap file)
  - ü Virtual address space not backed by "real" files also called "anonymous VM"

**n** Advantages
  - ü Uniform access for files and memory (just use pointers)
  - ü Less copying

**n** Drawbacks
  - ü Process has less control over data movement
    - § OS handles faults transparently
  - ü Does not generalize to streamed I/O (pipes, sockets, etc.)

# Page Replacement

**n** Prevent over-allocation of memory by modifying page-fault service routine to include page replacement

**n** Use *modify* (*dirty*) *bit* to reduce overhead of page transfers
  - ü Only modified pages are written to disk

**n** Page replacement completes separation between logical memory and physical memory
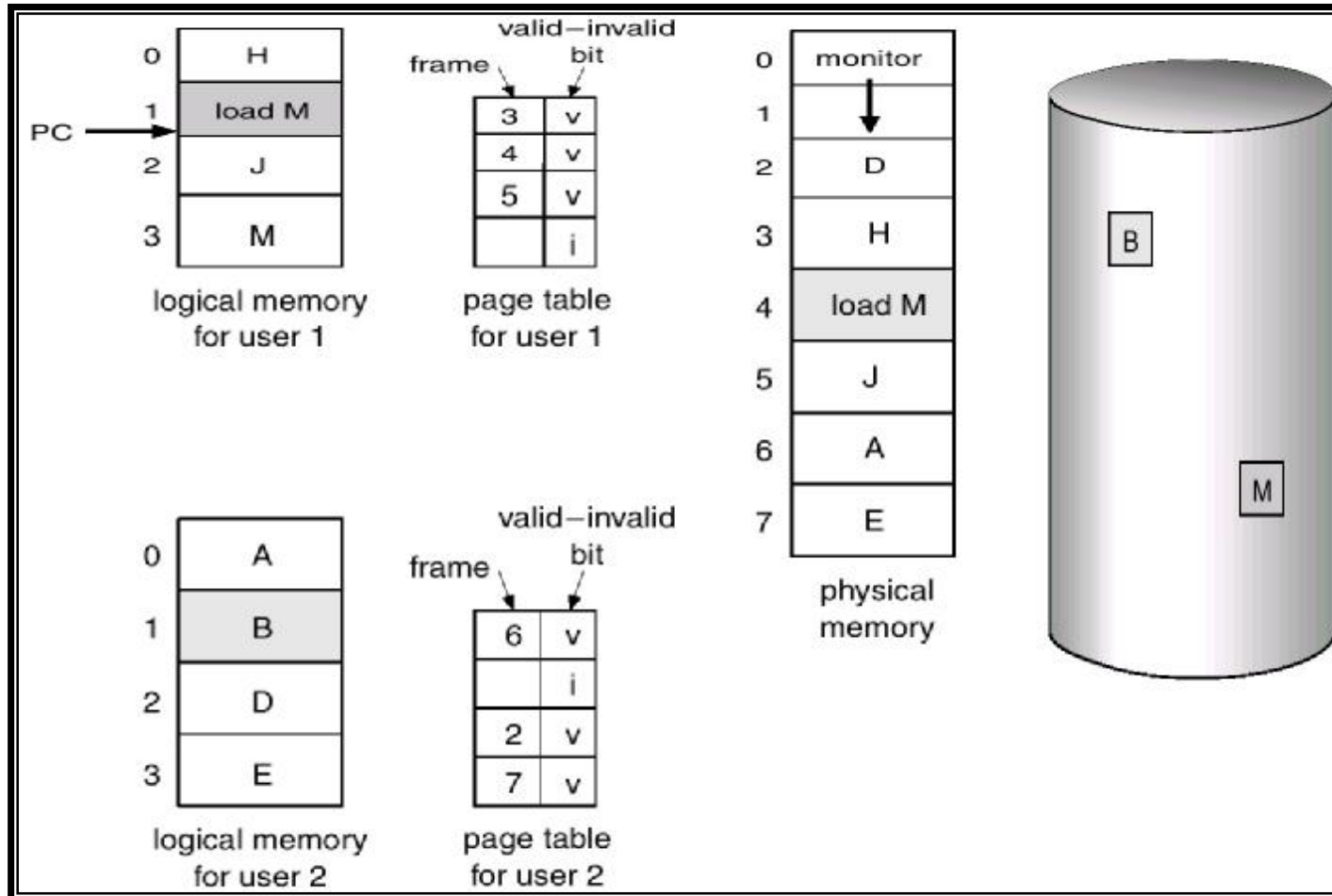  - ü Large virtual memory can be provided on a smaller physical memory

**n**  When a page fault occurs, the OS loads the faulted page from disk into a page frame of memory

**n**  At some point, the process has used all of the page frames it is allowed to use

**n**  When this happens, the OS must replace a page for each page faulted in

    **ü**  It must evict a page to free up a page frame

**n**  The page replacement algorithm determines how this is done

**n** Evicting the best page

- ü The goal of the replacement algorithm is to reduce the fault rate by selecting the best victim page to remove

- ü The best page to evict is the one never touched again
    - § as process will never again fault on it

- ü "Never" is a long time, so picking the page closest to "never" is the next best thing
    - § Belady's proof: Evicting the page that won't be used for the longest period of time minimizes the number of page faults
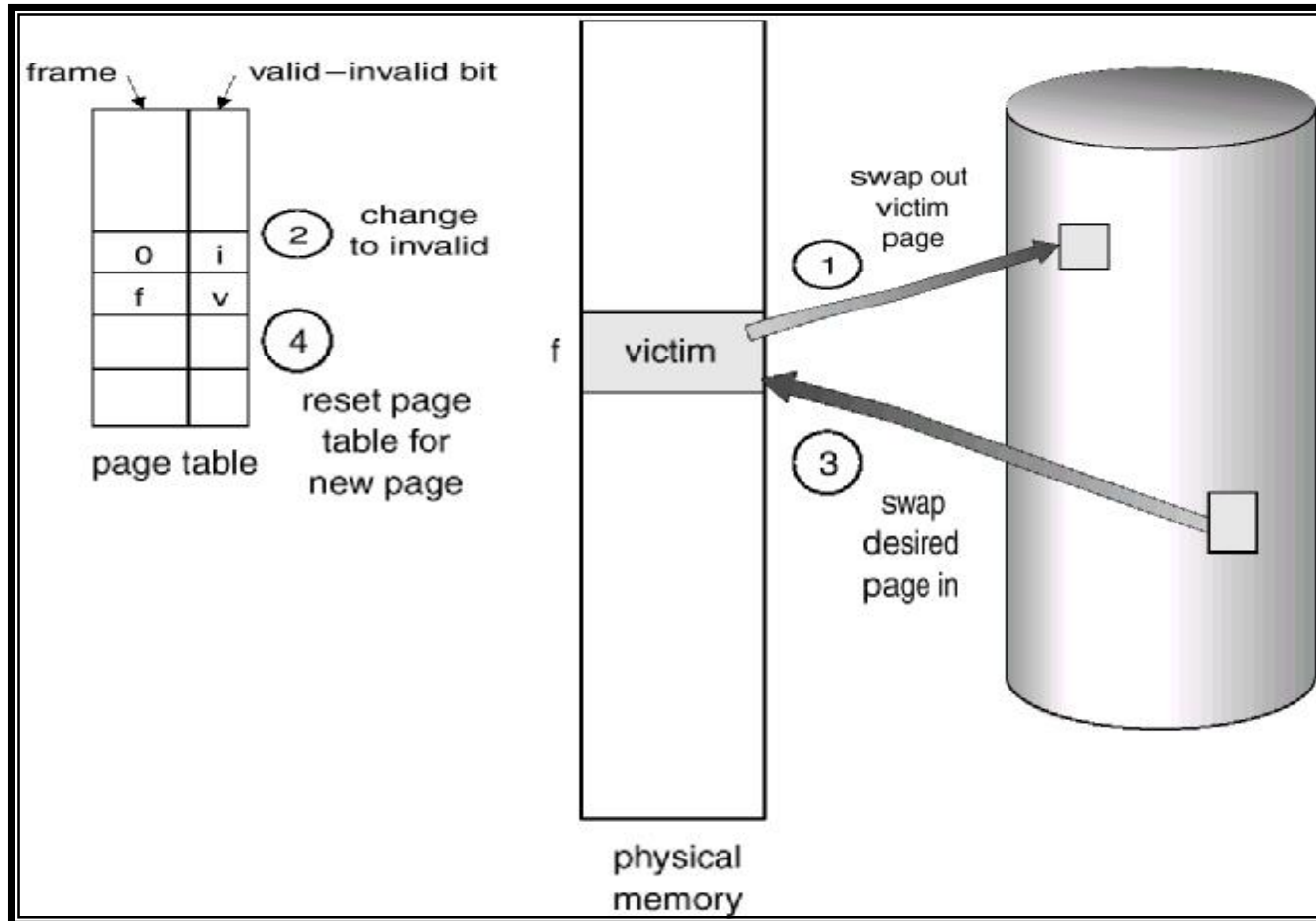
# Basic Page Replacement

1. Find the location of the desired page on disk

2. Find a free frame:
   - ü If there is a free frame, use it
   - ü If there is no free frame, use a page replacement algorithm to select a *victim* frame

3. Read the desired page into the (newly) free frame
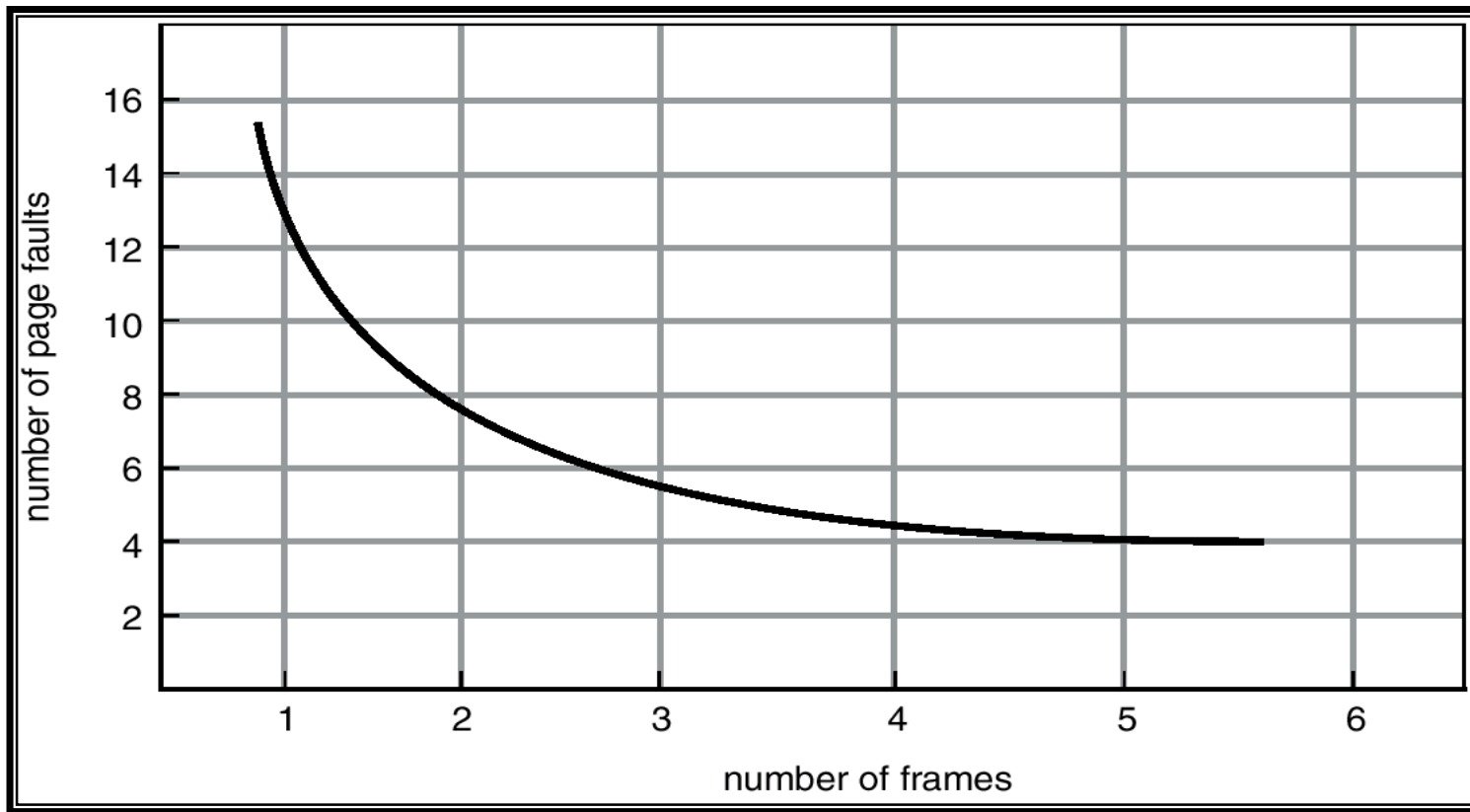   - ü Update the page and frame tables

4. Restart the process

# Page Replacement Algorithms

n  Want lowest page-fault rate

n  Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string

n  In all our examples, the reference string is

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

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

**n** Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

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

| 1 | 1 | 4 | 5 |
|---|---|---|---|
| 2 | 2 | 1 | 3 |
| 3 | 3 | 2 | 4 |

9 page faults

**n** 4 frames

| 1 | 1 | 5 | 4 |
|---|---|---|---|
| 2 | 2 | 1 | 5 |
| 3 | 3 | 2 | |
| 4 | 4 | 3 | |

10 page faults

**n** FIFO Replacement – Belady's Anomaly

   **ü** more frames ⇒ less page faults

**n** Obvious and simple to implement

   ü Maintain a list of pages in order they were paged in

   ü On replacement, evict the one brought in longest time ago

**n** Why might this be good?

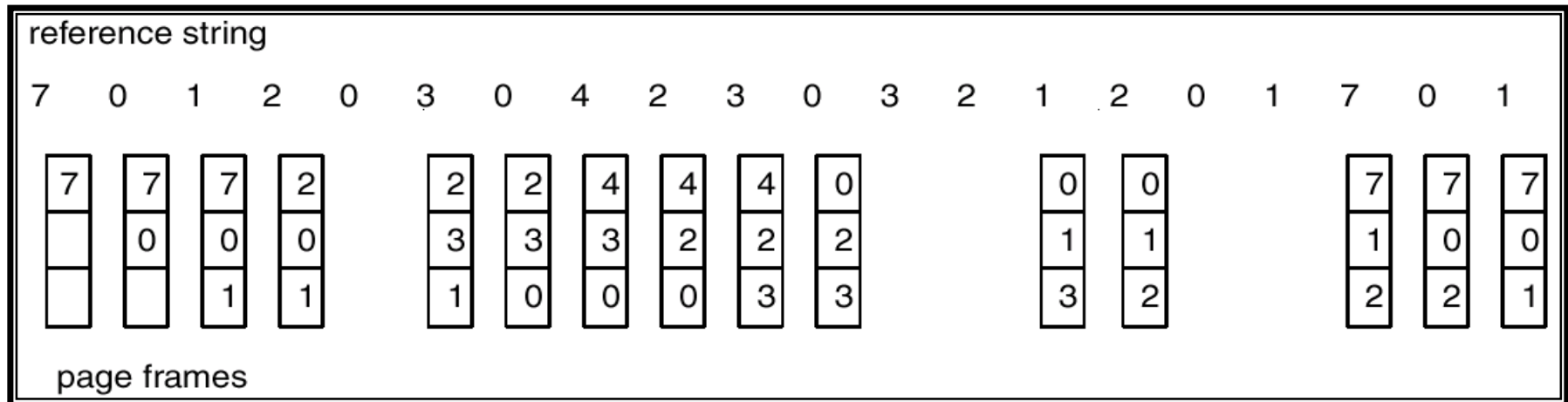   ü Maybe the one brought in the longest ago is not being used

**n** Why might this be bad?

   ü Maybe, it's not the case

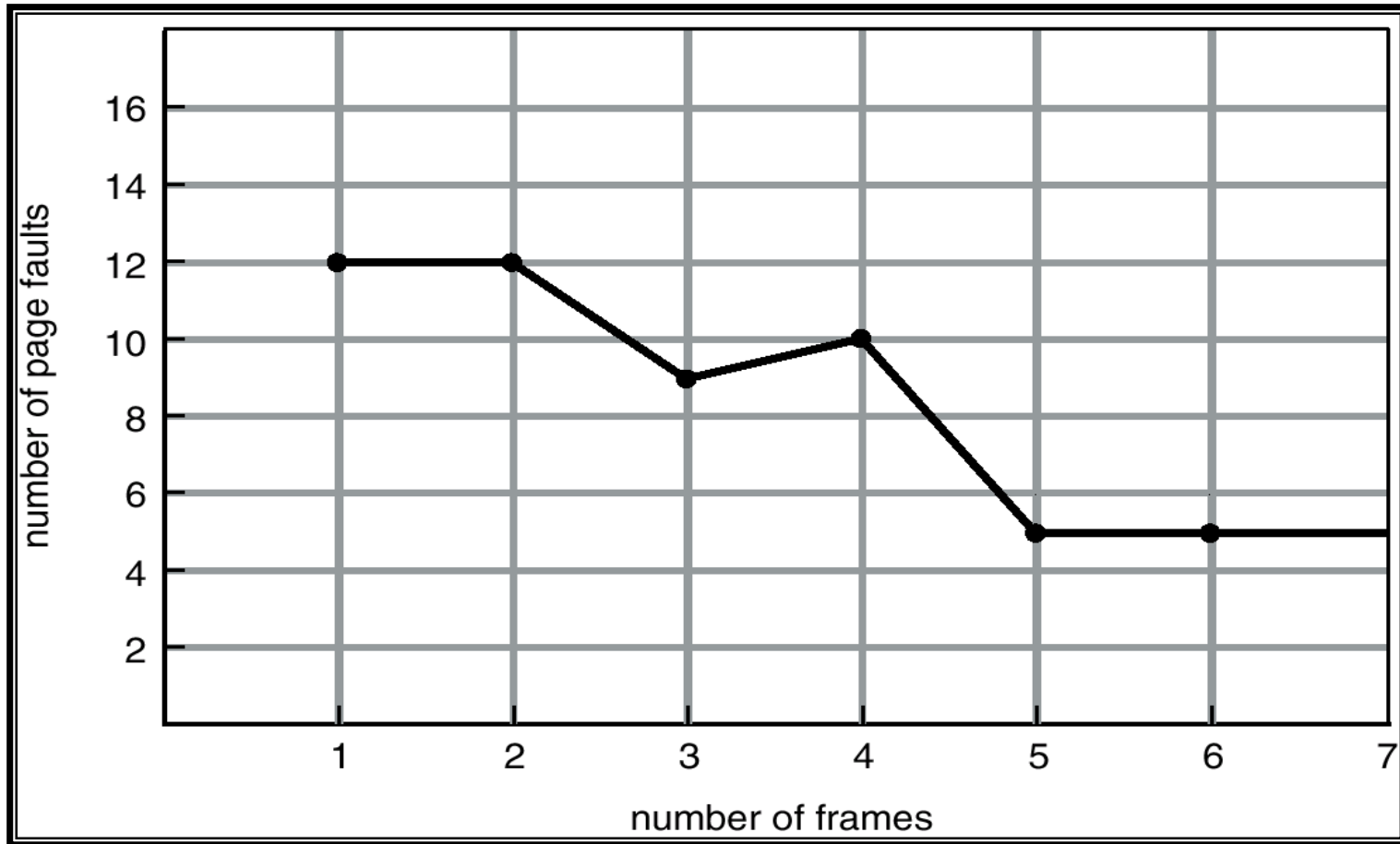   ü We don't have any information either way

**n** FIFO suffers from "Belady's Anomaly"

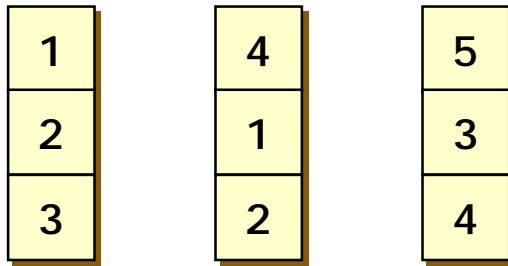   ü The fault rate might increase when the algorithm is given more memory

reference string

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

page frames

**n** Example: Belady's anomaly

   **ü** Reference string: 1,2,3,4,1,2,5,1,2,3,4,5

   **ü** 3 frames: 9 faults

| 1 |
|---|
| 2 |
| 3 |

| 4 |
|---|
| 1 |
| 2 |

| 5 |
|---|
| 3 |
| 4 |

   **ü** 4 frames: 10 faults

| 1 |
|---|
| 2 |
| 3 |
| 4 |

| 5 |
|---|
| 1 |
| 2 |
| 3 |

| 4 |
|---|
| 5 |

# Optimal Algorithm

n  Replace page that will not be used for longest period of time
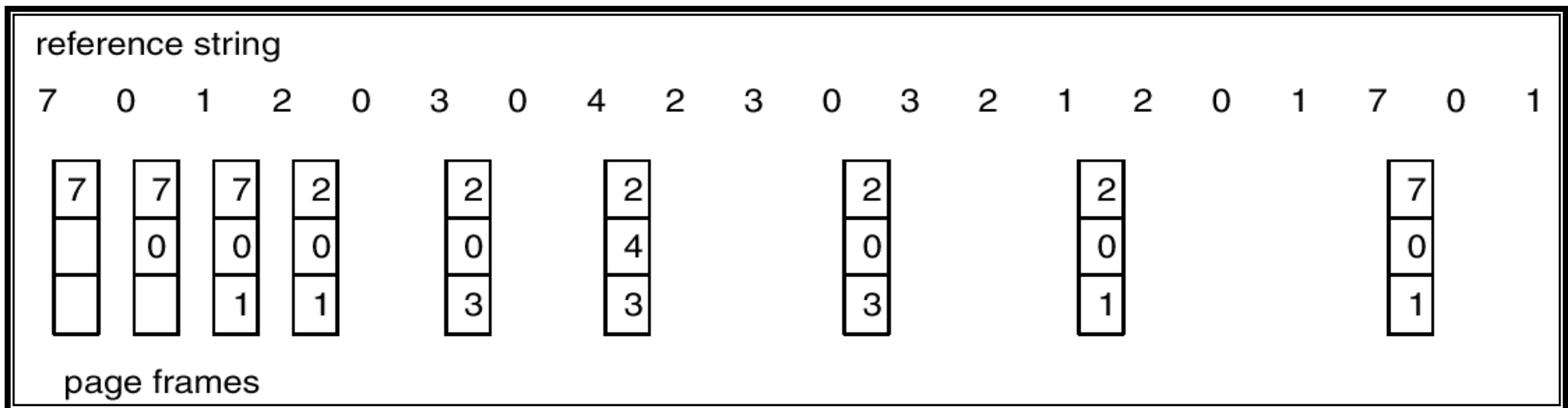
n  4 frames example

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

| | |
|---|---|
| 1 | 4 |
| 2 | |
| 3 | |
| 4 | 5 |

6 page faults

n  How do you know this?

n  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

# Least Recently Used (LRU) Algorithm

**n** Reference string: 1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

| 1 | 5 |
|---|---|
| 2 |   |
| 3 | 5   4 |
| 4 | 3 |

**n** Counter implementation

   ü Every page entry has a counter

   ü Every time page is referenced through this entry, copy the clock into the counter

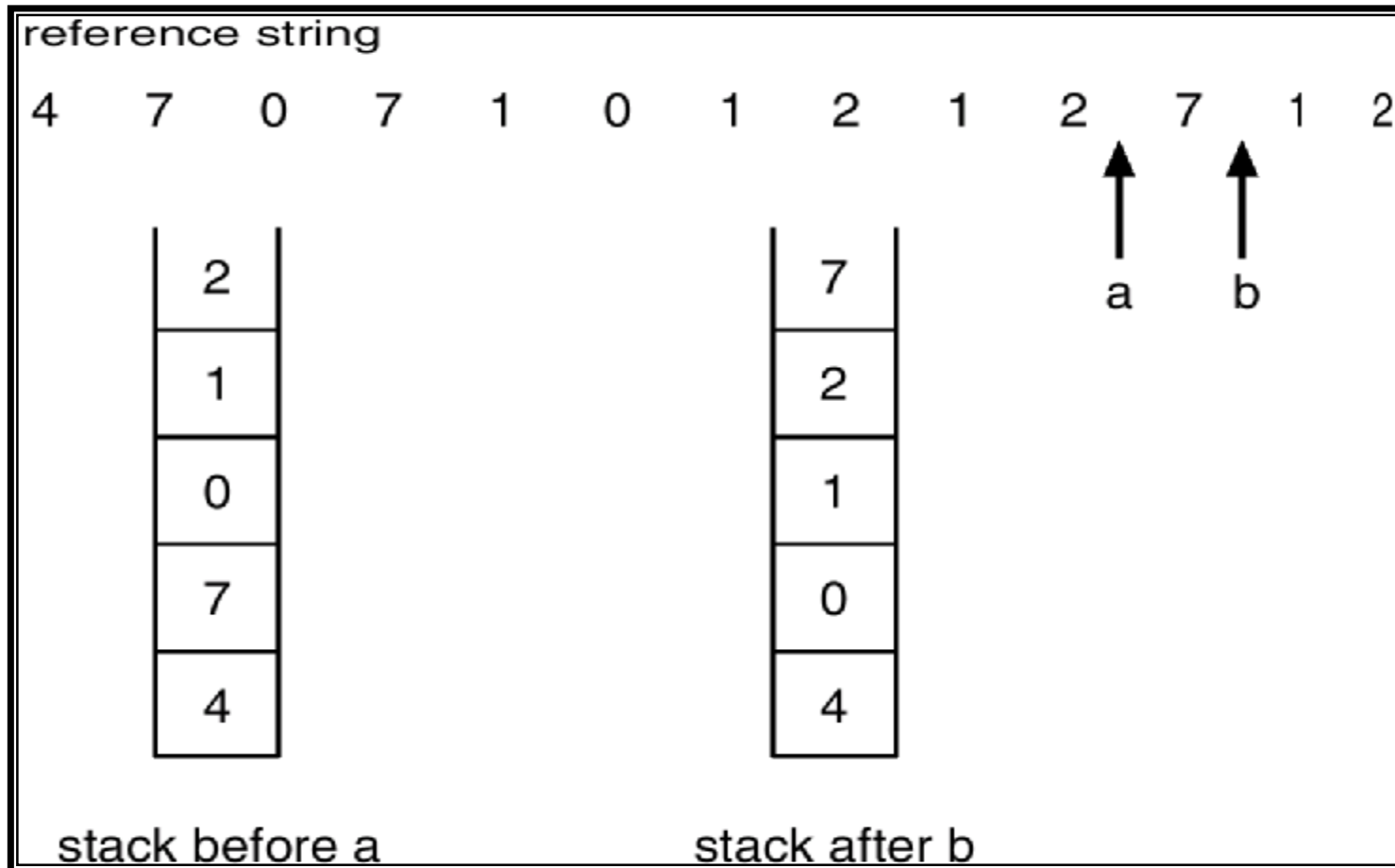   ü When a page needs to be changed, look at the counters to determine which are to change

reference string

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

page frames

# LRU Algorithm (Cont'd)

**n** Stack implementation – keep a stack of page numbers in a double link form:

    ü Page referenced:

        § move it to the top

        § requires 6 pointers to be changed

    ü No search for replacement

**n** LRU uses reference information to make a more informed replacement decision

- ü Idea: past experience gives us a guess of future behavior
- ü On replacement, evict the page that has not been used for the longest time in the past
- ü LRU looks at the past, Belady's wants to look at future

**n** Implementation

- ü To be perfect, need to timestamp every reference and put it in the PTE (or maintain a stack) – too expensive
- ü So, we need an approximation

# LRU Approximation Algorithms

**n** Reference bit

- ü With each page associate a bit, initially = 0
- ü When page is referenced bit set to 1
- ü Replace the one which is 0 (if one exists). We do not know the order, however

**n** Second chance

- ü Need reference bit
- ü Clock replacement
- ü If page to be replaced (in clock order) has reference bit = 1, then:
    - § set reference bit 0
    - § leave page in memory
    - § replace next page (in clock order), subject to same rules

**n** Many LRU approximations use the PTE reference (R) bit

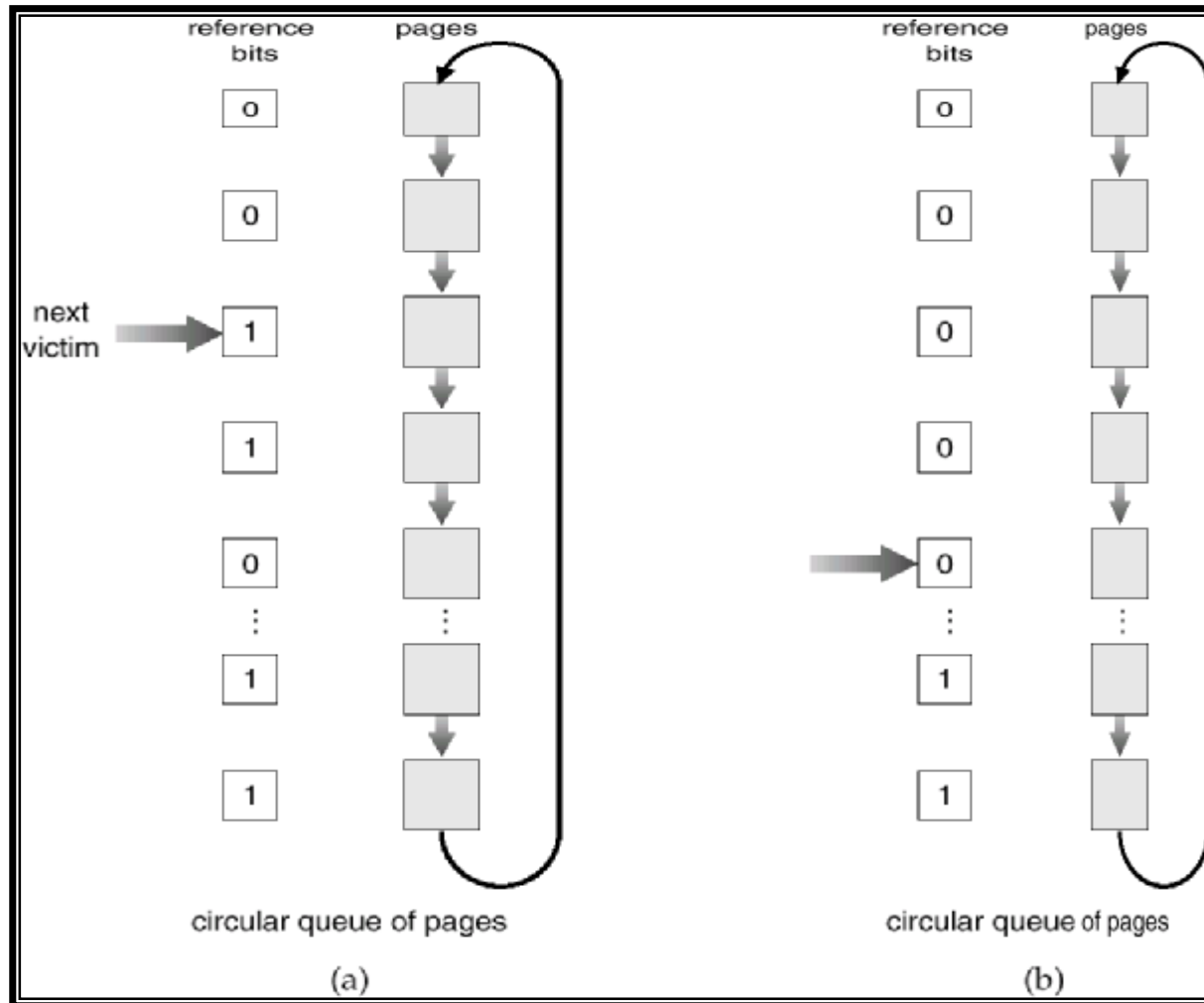    ü R bit is set whenever the page is referenced (read or written)

**n** Counter-based approach

    ü Keep a counter for each page

    ü At regular intervals, for every page, do:

        § If R = 0, increment the counter (hasn't been used)

        § If R = 1, zero the counter (has been used)

        § Zero the R bit

    ü The counter will contain the number of intervals since the last reference to the page

    ü The page with largest counter is the least recently used

**n** Some architectures don't have a reference bit

    ü Can simulate reference bit using the valid bit to induce faults

# Second Chance or LRU Clock

n  FIFO with giving a second chance to a recently referenced page

n  Arrange all of physical page frames in a big circle (clock)

n  A clock hand is used to select a good LRU candidate
  - ü Sweep through the pages in circular order like a clock
  - ü If the R bit is off, it hasn't been used recently and we have a victim
  - ü If the R bit is on, turn it off and go to next page

n  Arm moves quickly when pages are needed
  - ü Low overhead if we have plenty of memory
  - ü If memory is large, "accuracy" of information degrades
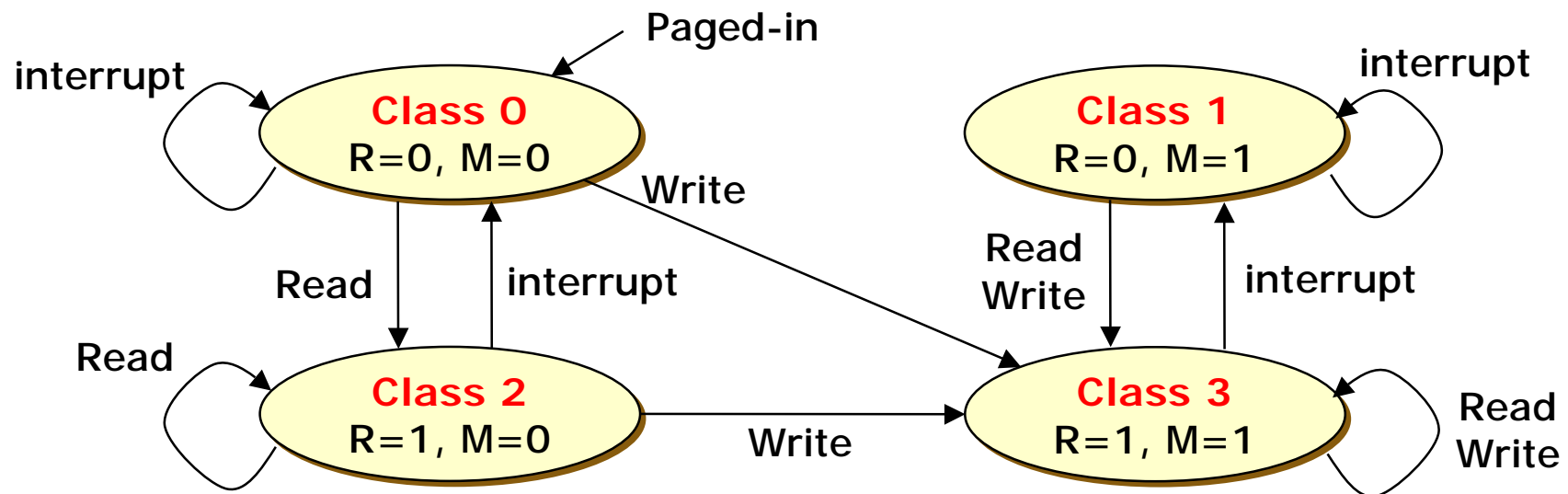
**n** NRU or enhanced second chance

ü Use R (reference) and M (modify) bits

§ Periodically, (e.g., on each clock interrupt), R is cleared, to distinguish pages that have not been referenced recently from those that have been

**n** Algorithm

- ü Removes a page at random from the lowest numbered nonempty class
- ü It is better to remove a modified page that has not been referenced in at least one clock tick than a clean page that is in heavy use

**n** Advantages

- ü Easy to understand
- ü Moderately efficient to implement
- ü Gives a performance that, while certainly not optimal, may be adequate

# Counting Algorithms

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

**n** LFU Algorithm: replaces page with smallest count

**n** MFU Algorithm: based on the argument that the page with the smallest count was probably just brought in and has yet to be used

**n** Counting-based page replacement

ü A software counter is associated with each page

ü At each clock interrupt, for each page, the R bit is added to the counter

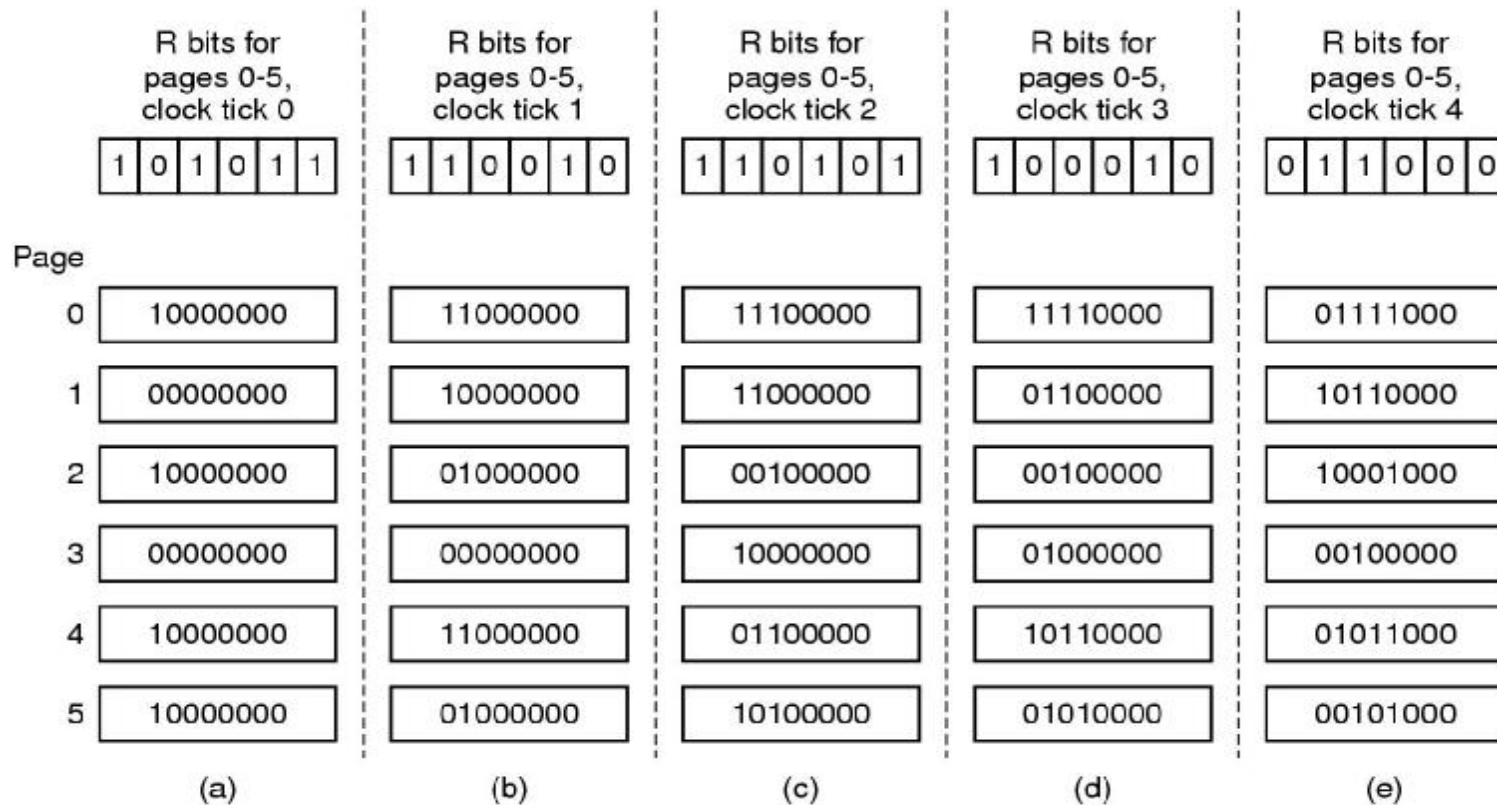§ The counters denote how often each page has been referenced

**n** Least Frequently Used (LFU)

ü The page with the smallest count will be replaced

ü Cf) Most frequently used (MFU) page replacement

§ The page with the largest count will be replaced

§ Based on the argument that the page with the smallest count was probably just brought in and has yet to be used

ü It never forgets anything

§ A page may be heavily used during the initial phase of a process, but then is never used again

**n Aging**

ü The counters are shifted right by 1 bit before the R bit is added to the leftmost

| | R bits for pages 0-5, clock tick 0 | R bits for pages 0-5, clock tick 1 | R bits for pages 0-5, clock tick 2 | R bits for pages 0-5, clock tick 3 | R bits for pages 0-5, clock tick 4 |
|---|---|---|---|---|---|
| | 1 0 1 0 1 1 | 1 1 0 0 1 0 | 1 1 0 1 0 1 | 1 0 0 0 1 0 | 0 1 1 0 0 0 |
| Page | | | | | |
| 0 | 10000000 | 11000000 | 11100000 | 11110000 | 01111000 |
| 1 | 00000000 | 10000000 | 11000000 | 01100000 | 10110000 |
| 2 | 10000000 | 01000000 | 00100000 | 00100000 | 10001000 |
| 3 | 00000000 | 00000000 | 10000000 | 01000000 | 00100000 |
| 4 | 10000000 | 11000000 | 01100000 | 10110000 | 01011000 |
| 5 | 10000000 | 01000000 | 10100000 | 01010000 | 00101000 |
| | (a) | (b) | (c) | (d) | (e) |

# Allocation of Frames

**n** Each process needs **minimum** number of pages

**n** Example: IBM 370 – 6 pages to handle SS MOVE instruction:

    ü instruction is 6 bytes, might span 2 pages

    ü 2 pages to handle **from**

    ü 2 pages to handle **to**

**n** Two major allocation schemes

    ü fixed allocation

    ü priority allocation

# Fixed Allocation

n   Equal allocation – e.g., if 100 frames and 5 processes, give each 20 pages

n   Proportional allocation – Allocate according to the size of process

- $s_i = $ size of process $p_i$
- $S = \sum s_i$
- $m = $ total number of frames
- $a_i = $ allocation for $p_i = \dfrac{s_i}{S} \times m$

$$m = 64$$

$$s_1 = 10$$

$$s_2 = 127$$

$$a_1 = \frac{10}{137} \times 64 \approx 5$$

$$a_2 = \frac{127}{137} \times 64 \approx 59$$

# Priority Allocation

**n**  Use a proportional allocation scheme using priorities rather than size

**n**  If process $P_i$ generates a page fault,
- ü  select for replacement one of its frames
- ü  select for replacement a frame from a process with lower priority number

# Global vs. Local Allocation

**n  Global** replacement

    ü  Process selects a replacement frame from the set of all frames

    ü  One process can take a frame from another


**n  Local** replacement

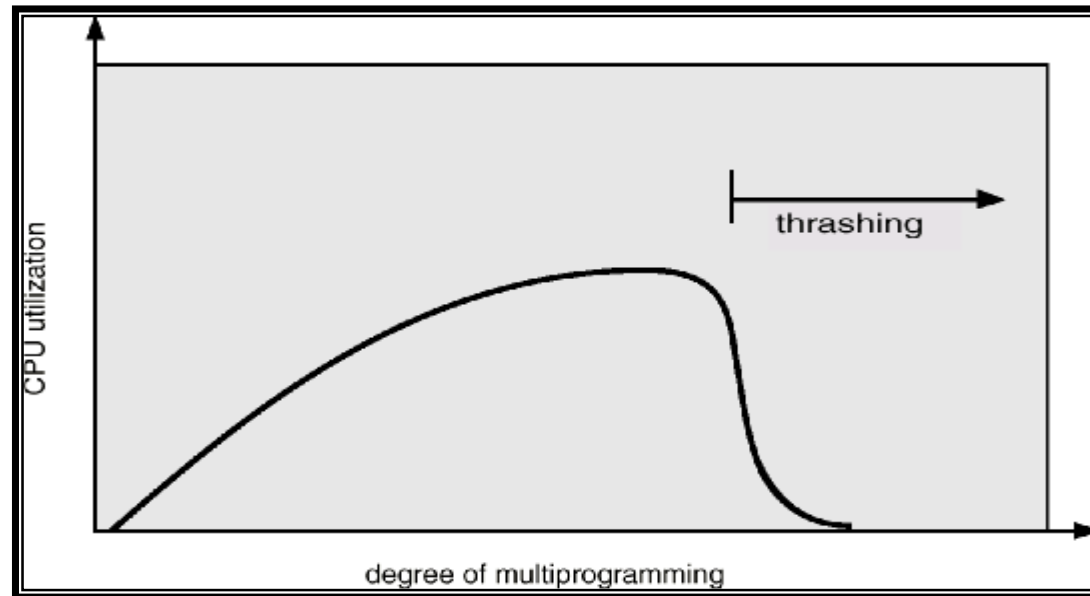    ü  Each process selects from only its own set of allocated frames

# Thrashing

n  If a process does not have "enough" pages, the page-fault rate is very high

n  This leads to:
  ü  Low CPU utilization
  ü  Operating system thinks that it needs to increase the degree of multiprogramming
  ü  Another process added to the system

n  **Thrashing** $\equiv$ a process is busy swapping pages in and out
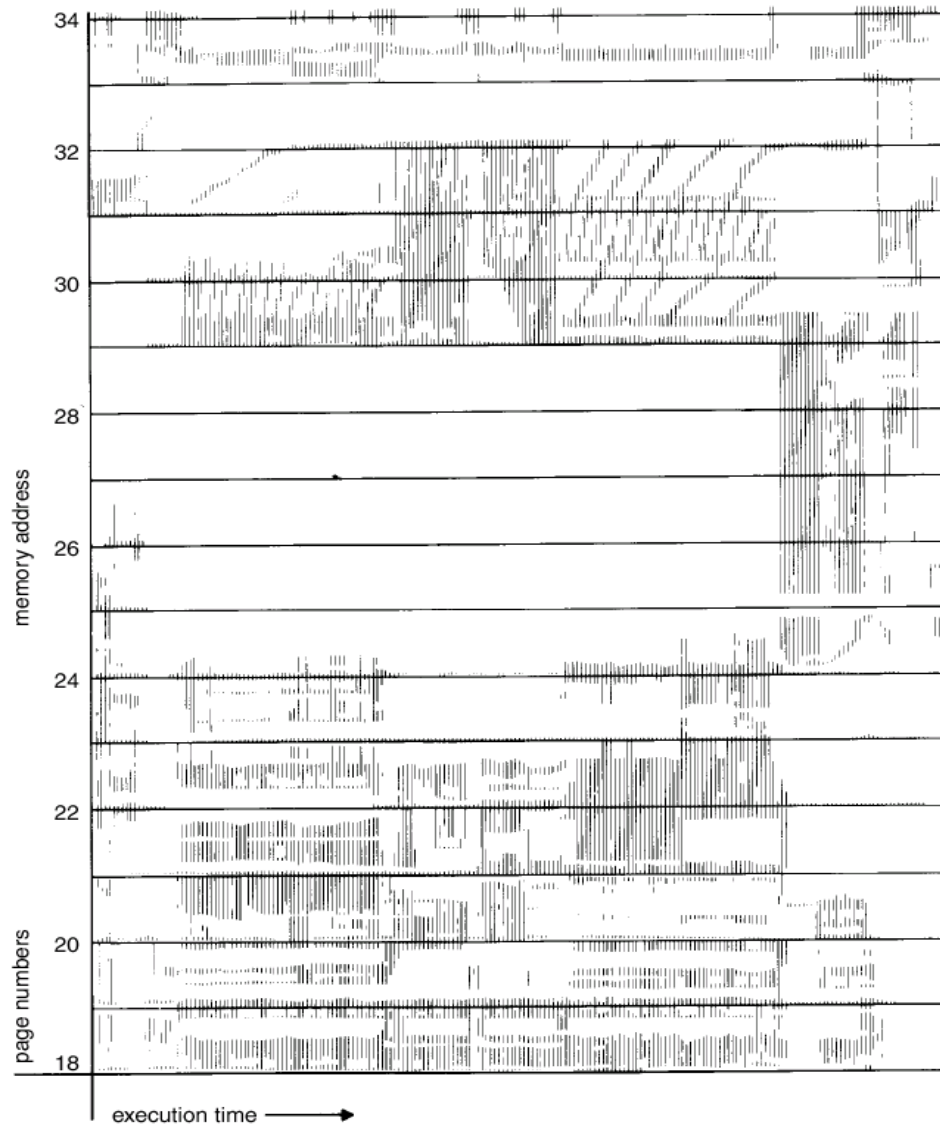
# Thrashing



- **n** Why does paging work? **à** Locality model
    - **ü** Process migrates from one locality to another
    - **ü** Localities may overlap


- **n** Why does thrashing occur?
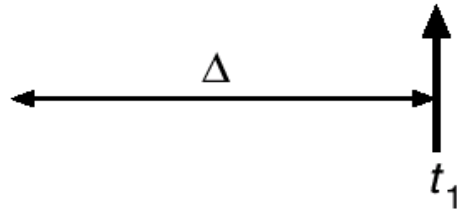    - **ü** $\Sigma$ size of locality > total memory size

# Working-Set Model

n $\Delta \equiv$ working-set window $\equiv$ a fixed number of page references
- ü Example: 10,000 instruction

n $WSS_i$ (Working Set Size of Process $P_i$) =
  total number of pages referenced in the most recent $\Delta$ (varies in time)
- ü if $\Delta$ too small will not encompass entire locality
- ü if $\Delta$ too large will encompass several localities
- ü if $\Delta = \infty \Rightarrow$ will encompass entire program

n $D = \Sigma\ WSS_i \equiv$ total demand frames

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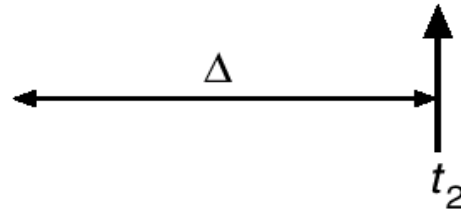
n Policy if $D > m$, then suspend one of the processes

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

$\Delta$       $t_1$       $\Delta$       $t_2$
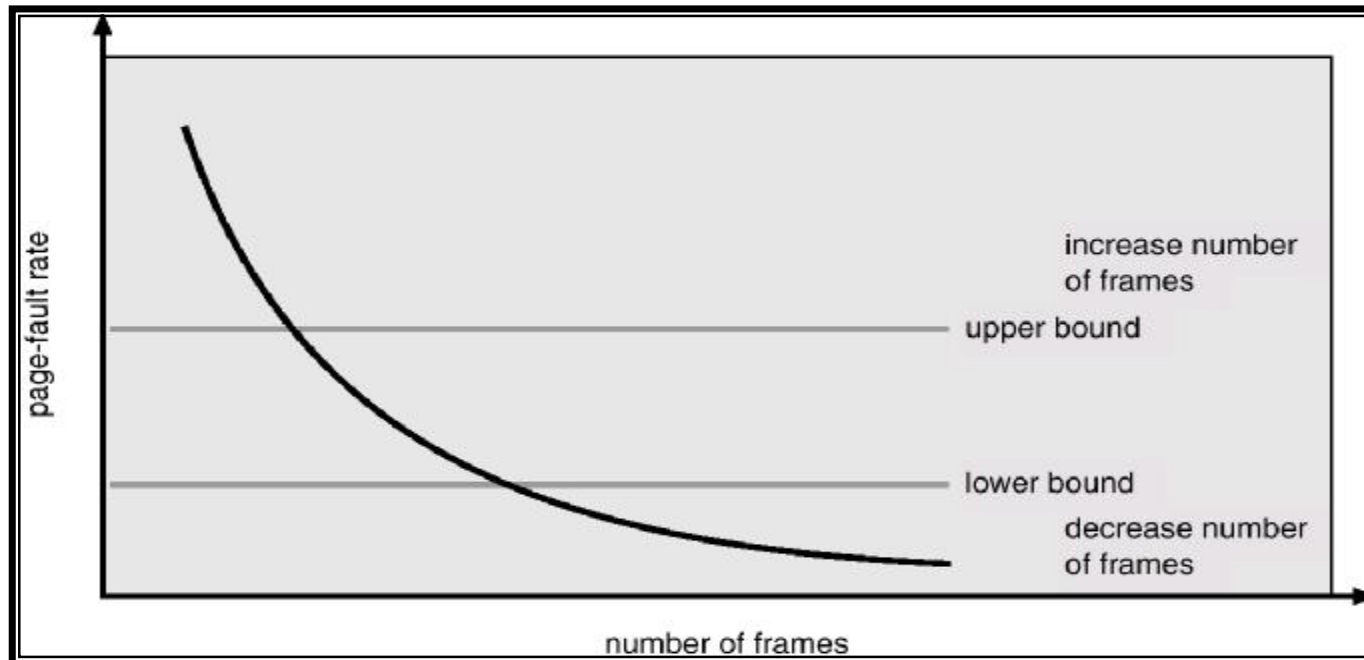
$WS(t_1) = \{1,2,5,6,7\}$       $WS(t_2) = \{3,4\}$

# Keeping Track of the Working Set

n  Approximate with interval timer + a reference bit

n  Example: $\Delta$ = 10,000
  ü  Timer interrupts after every 5000 time units
  ü  Keep in memory 2 bits for each page
  ü  Whenever a timer interrupts copy and sets the values of all reference bits to 0
  ü  If one of the bits in memory = 1 $\Rightarrow$ page in working set

n  Why is this not completely accurate?

n  Improvement à 10 bits and interrupt every 1000 time units

# Page-Fault Frequency Scheme



**n** Establish "acceptable" page-fault rate

    ù If actual rate too low, process loses frame

    ù If actual rate too high, process gains frame

# Other Considerations

**n** Prepaging

**n** Page size selection
- ü Fragmentation
- ü Table size
- ü I/O overhead
- ü Locality

# Other Considerations (Cont'd)

**n** TLB Reach

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

# Increasing the Size of the TLB

**n** Increase the Page Size

    **ü** This may lead to an increase in fragmentation as not all applications require a large page size

**n** Provide Multiple Page Sizes

    **ü** This allows applications that require larger page sizes the opportunity to use them without an increase in fragmentation

# Other Considerations (Cont'd)

**n** Program structure

    ü **int A[][] = new int[1024][1024];**

    ü Each row is stored in one page

    ü Program 1

```
for (j = 0; j < A.length; j++)
    for (i = 0; i < A.length; i++)
        A[i,j] = 0;
```

    1024 x 1024 page faults

    ü Program 2

```
for (i = 0; i < A.length; i++)
    for (j = 0; j < A.length; j++)
        A[i,j] = 0;
```

    1024 page faults
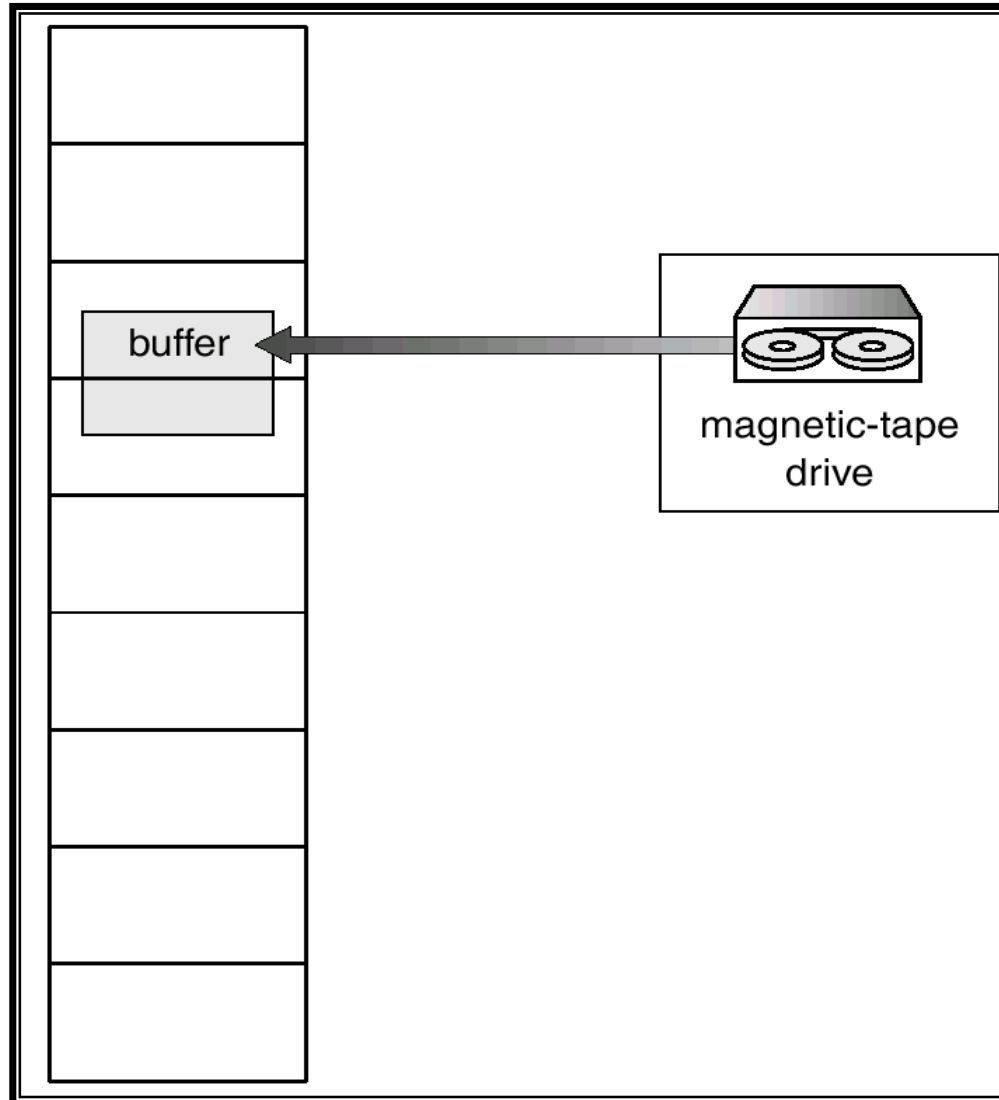
# Other Considerations (Cont'd)

**n** I/O Interlock

  ü Pages must sometimes be locked into memory

**n** Consider I/O

  ü Pages that are used for copying a file from a device must be locked from being selected for eviction by a page replacement algorithm

# Operating System Examples

**n** Windows NT

**n** Solaris 2

# Windows NT

**n** Uses demand paging with **clustering**
  - ü Clustering brings in pages surrounding the faulting page

**n** Processes are assigned **working set minimum** and **working set maximum**
  - ü Working set minimum is the minimum number of pages the process is guaranteed to have in memory
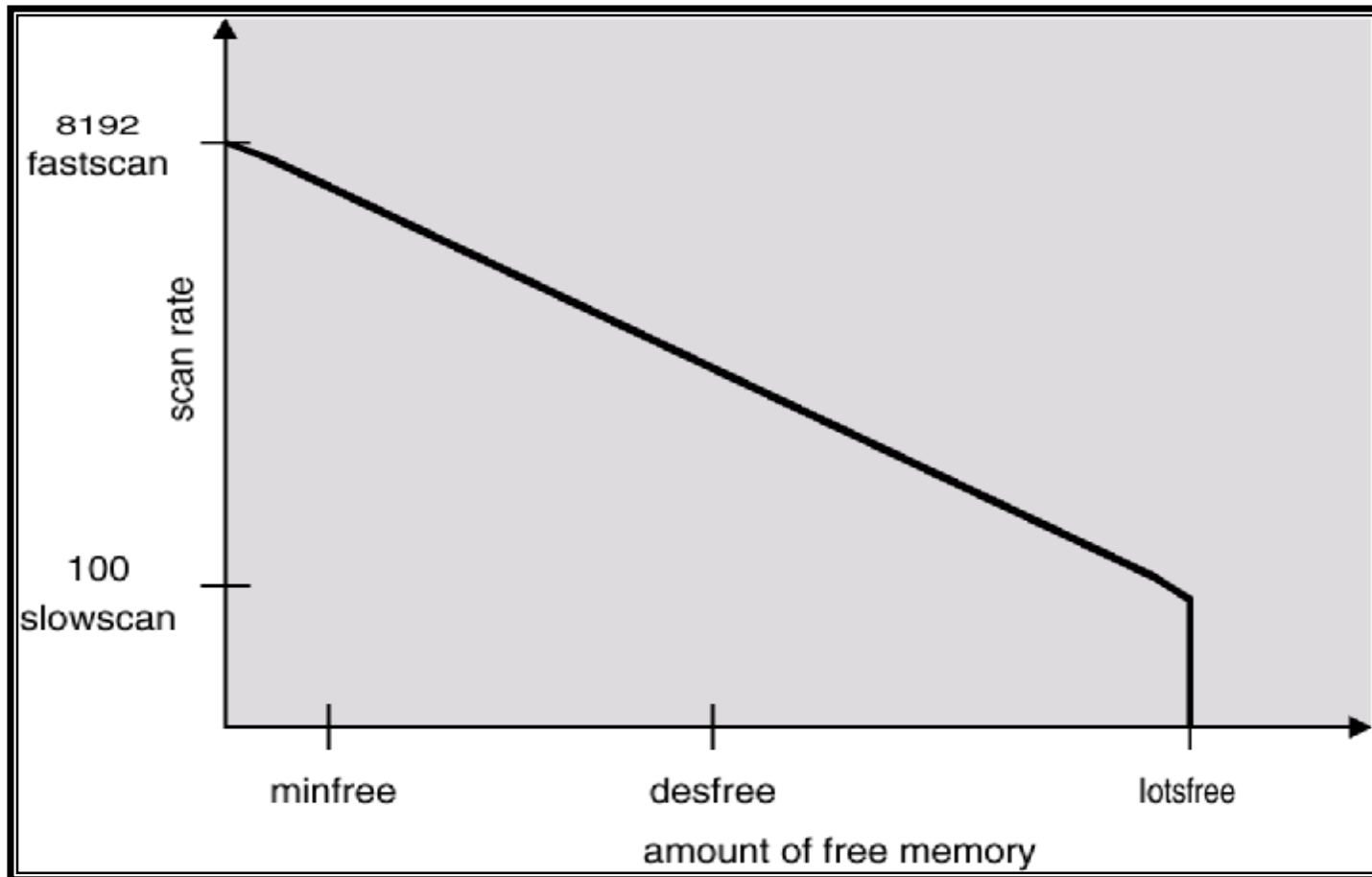  - ü A process may be assigned as many pages up to its working set maximum

**n** 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

**n** Working set trimming removes pages from processes that have pages in excess of their working set minimum

# Solaris 2

**n** Maintains a list of free pages to assign faulting processes

**n** **Lotsfree** – threshold parameter to begin paging

**n** Paging is peformed by *pageout* process

**n** Pageout scans pages using modified clock algorithm
  - ü Two-handed-clock algorithm (similar to the second-chance algorithm)
  - ü *handspread*

**n** **Scanrate** is the rate at which pages are scanned
  - ü This ranged from **slowscan** to **fastscan**

**n** Pageout is called more frequently depending upon the amount of free memory available

**n** Advantages

ü Separates user's logical memory from physical memory

§ Abstracts main memory into an extremely large, uniform array of storage

§ Frees programmers from the concerns of memory-storage limitations

ü Allows the execution of processes that may not be completely in memory

§ Programs can be larger than physical memory

§ More programs could be run at the same time

§ Less I/O would be needed to load or swap each user program into memory

ü Allows processes to easily share files and address spaces

ü Provides an efficient mechanism for protection and process creation

**n** Disadvantages

ü Performance!!!

§ In terms of time and space

n Optimizations
- ü Managing page tables (space)
- ü Efficient Translation (TLBs) (time)
- ü Demand paging (space)

n Advanced functionality
- ü Sharing memory
- ü Copy on write
- ü Mapped files