# Chapter 9:  Virtual Memory

# Chapter 9: Virtual Memory

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

# 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

# Background

- Preceding sections talked about how to avoid memory fragmentation by breaking process memory requirements down into smaller bites ( pages ), and storing the pages non-contiguously in memory.

- However the entire process still had to be stored in memory somewhere.

- In practice, most real processes do not need all their pages, or at least not all at once, for several reasons:
  - Error handling code is not needed unless that specific error occurs, some of which are quite rare.
  - Arrays are often over-sized for worst-case scenarios, and only a small fraction of the arrays are actually used in practice.
  - Certain features of certain programs are rarely used, such as the routine to balance the federal budget.

# Background

- The ability to load only the portions of processes that were actually needed has several benefits:

- Programs could be written for a much larger address space ( virtual memory space ) than physically exists on the computer.

- Because each process is only using a fraction of their total address space, there is more memory left for other programs

  - improving CPU utilization and system throughput.

- Less I/O is needed for swapping processes in and out of RAM, speeding things up.
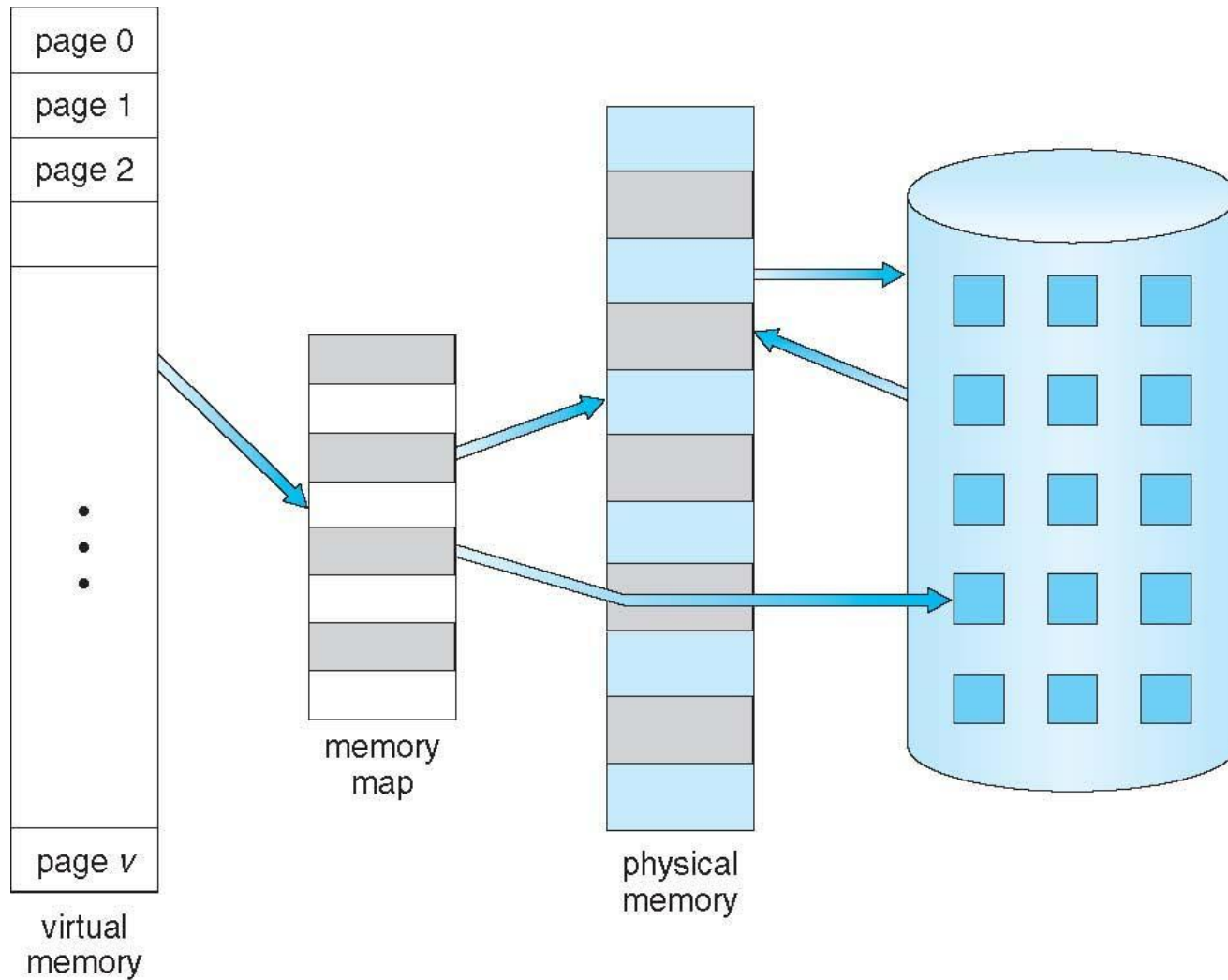
# Background (Cont.)

- **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
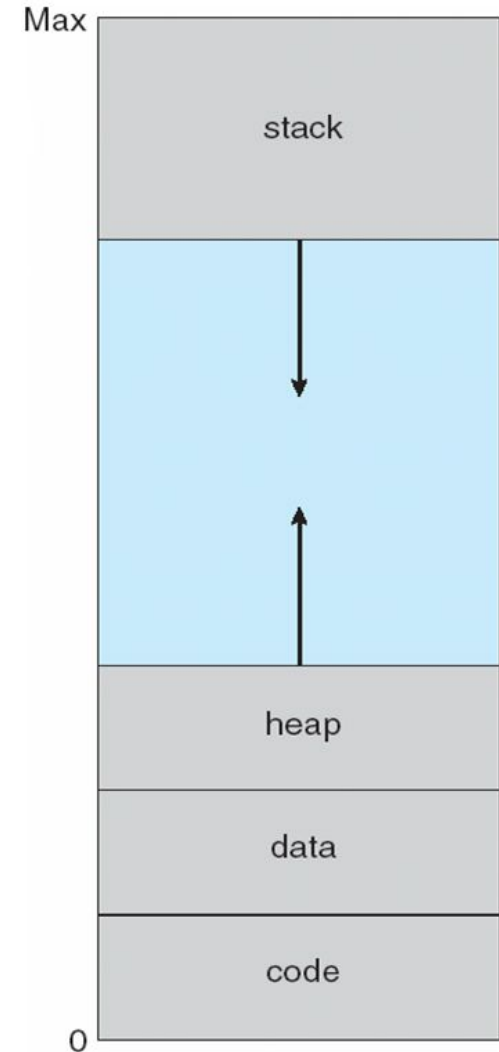  - Less I/O needed to load or swap processes

# 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

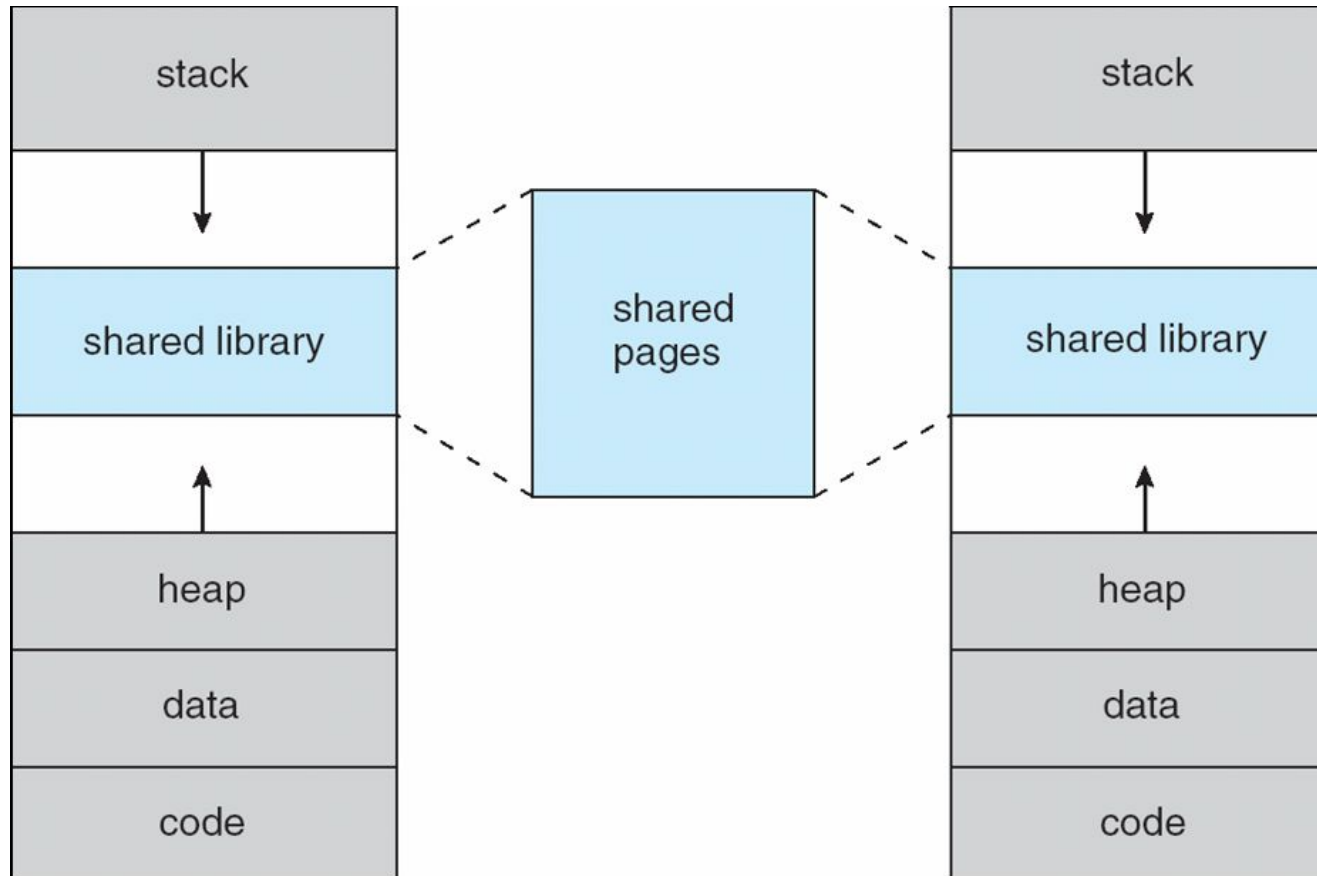# Virtual Memory That is Larger Than Physical Memory

# 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
    - 4  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
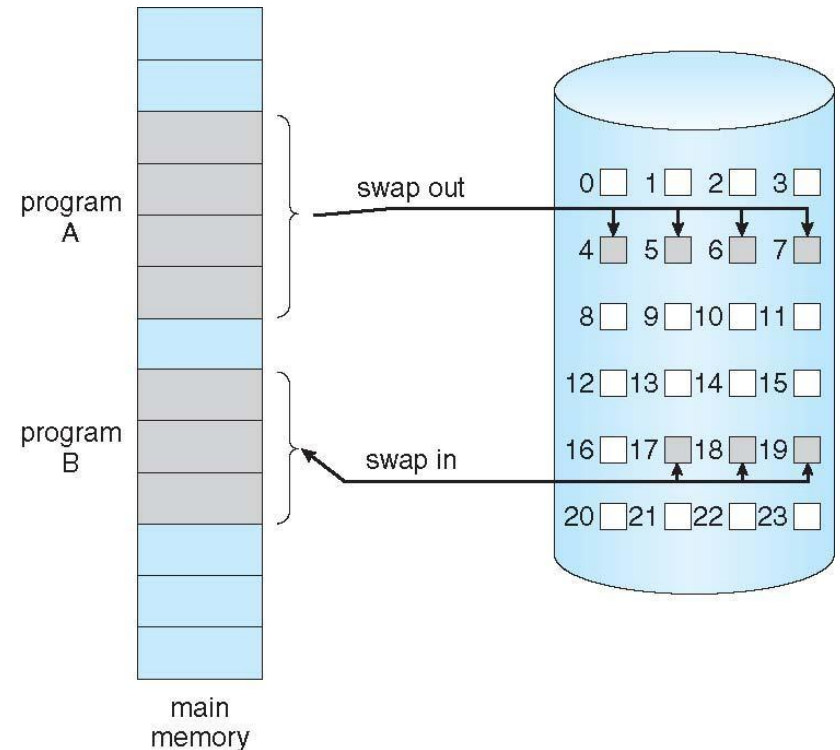- Pages can be shared during `fork()`, speeding process creation

# 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

# Basic Concepts

- The basic idea behind paging is that when a process is swapped in, the pager only loads into memory those pages that it expects the process to need

- Pages that are not loaded into memory are marked as invalid in the page table, using the invalid bit.

- If the process only ever accesses pages that are loaded in memory ( memory resident pages ), then the process runs exactly as if all the pages were loaded in to memory.

# 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:

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

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

page table

# Page Table When Some Pages Are Not in Main Memory



logical memory

frame    valid–invalid bit

page table

physical memory

# Page Fault

- On the other hand, if a page is needed that was not originally loaded up, then a page fault trap is generated, which must be handled in a series of steps:
  - The memory address requested is first checked, to make sure it was a valid memory request.
  - If the reference was invalid, the process is terminated. Otherwise, the page must be paged in.
  - A free frame is located, possibly from a free-frame list.
  - A disk operation is scheduled to bring in the necessary page from disk.
  - When the I/O operation is complete, the process's page table is updated with the new frame number, and the invalid bit is changed to indicate that this is now a valid page reference.
  - The instruction that caused the page fault must now be restarted from the beginning,
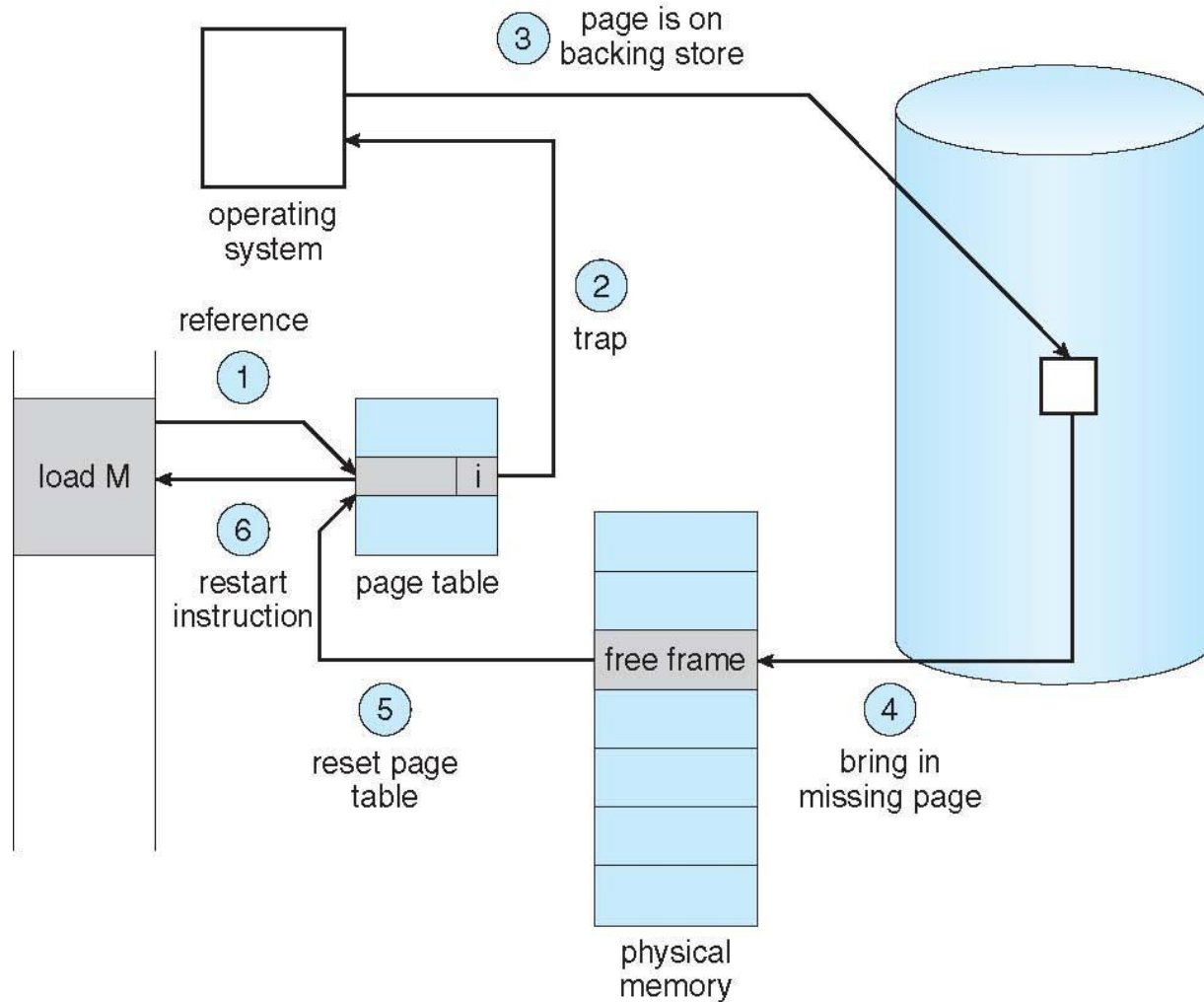
# Page Fault

- 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 of process, non-memory-resident -> page fault
  - And for every other process pages on first access
  - **Pure demand paging**
- Actually, a given instruction could access multiple pages -> multiple page faults
  - Consider fetch and decode of instruction which adds 2 numbers from memory and stores result back to memory
  - Pain decreased because of **locality of reference**
- Hardware support needed for demand paging
  - Page table with valid / invalid bit
  - Secondary memory (swap device with **swap space**)
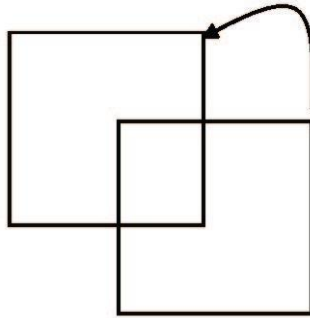  - Instruction restart

# Instruction Restart

- A crucial part of the process is that the instruction must be restarted from scratch once the desired page has been made available in memory.

- For most simple instructions this is not a major difficulty.

- However there are some architectures that allow a single instruction to modify a fairly large block of data, ( which may span a page boundary ), and if some of the data gets modified before the page fault occurs, this could cause problems.

- One solution is to access both ends of the block before executing the instruction, guaranteeing that the necessary pages get paged in before the instruction begins.

# 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?

# 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.)

- Three major activities
  - Service the interrupt – careful coding means just several hundred instructions needed
  - Read the page – lots of time
  - Restart the process – again just a small amount of time
- 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 )

# Demand Paging Example

- suppose that a normal memory access requires 200 nanoseconds, and that servicing a page fault takes 8 milliseconds. ( 8,000,000 nanoseconds, or 40,000 times a normal memory access. ) With a *page fault rate* of p, ( on a scale from 0 to 1 ), the effective access time is now:

- Memory access time = 200 nanoseconds

- Average page-fault service time = 8 milliseconds

- EAT = $(1 - p)$ x 200 + p (8 milliseconds)

    = $(1 - p)$ x 200 + p x 8,000,000

    = 200 + p x 7,999,800

- If one access out of 1,000 causes a page fault, then

    EAT = 8.2 microseconds.

  This is a slowdown by a factor of 40!!

- If want performance degradation < 10 percent
    - 220 > 200 + 7,999,800 x p
      20 > 7,999,800 x p
    - p < .0000025
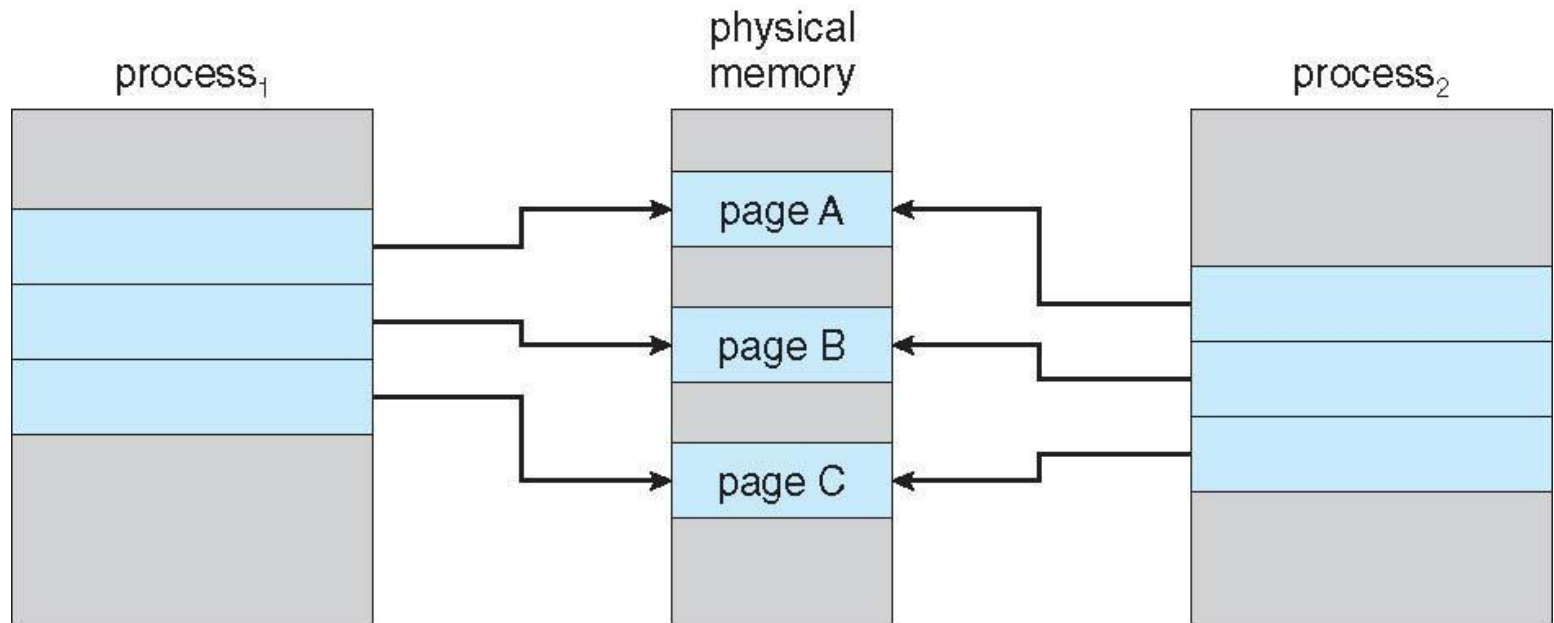    - < one page fault in every 400,000 memory accesses

# Demand Paging Optimizations

- Swap space I/O faster than file system I/O even if on the same device
  - Swap allocated in larger chunks, less management needed than file system

- Copy entire process image to swap space at process load time
  - Then page in and out of swap space
  - Used in older BSD Unix

- Demand page in from program binary on disk, but discard rather than paging out when freeing frame
  - Used in Solaris and current BSD
  - Still need to write to swap space
    - Pages not associated with a file (like stack and heap) – **anonymous memory**
    - Pages modified in memory but not yet written back to the file system

- Mobile systems
  - Typically don't support swapping
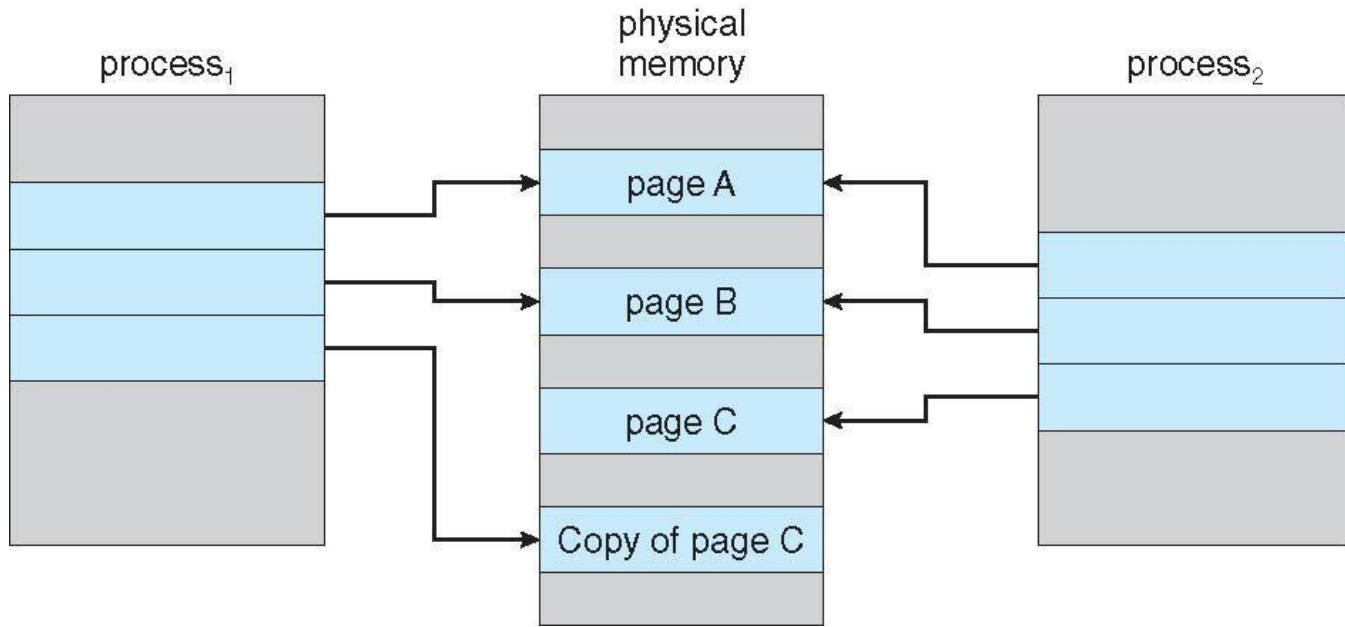  - Instead, demand page from file system and reclaim read-only pages (such as code)

# Copy-on-Write

- **Copy-on-Write** (COW) allows both parent and child processes to initially *share* the same pages in memory
    - If either process modifies a shared page, only then is the page copied
- COW allows more efficient process creation as only modified pages are copied
- In general, free pages are allocated from a **pool** of **zero-fill-on-demand** pages
    - Pool should always have free frames for fast demand page execution
        - Don't want to have to free a frame as well as other processing on page fault
    - Why zero-out a page before allocating it?
- `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?

- Adjust the memory used by I/O buffering, etc., to free up some frames for user processes.

- The decision of how to allocate memory for I/O versus user processes is a complex one, yielding different policies on different systems.
  - Some allocate a fixed amount for I/O
  - Others let the I/O system contend for memory along with everything else.

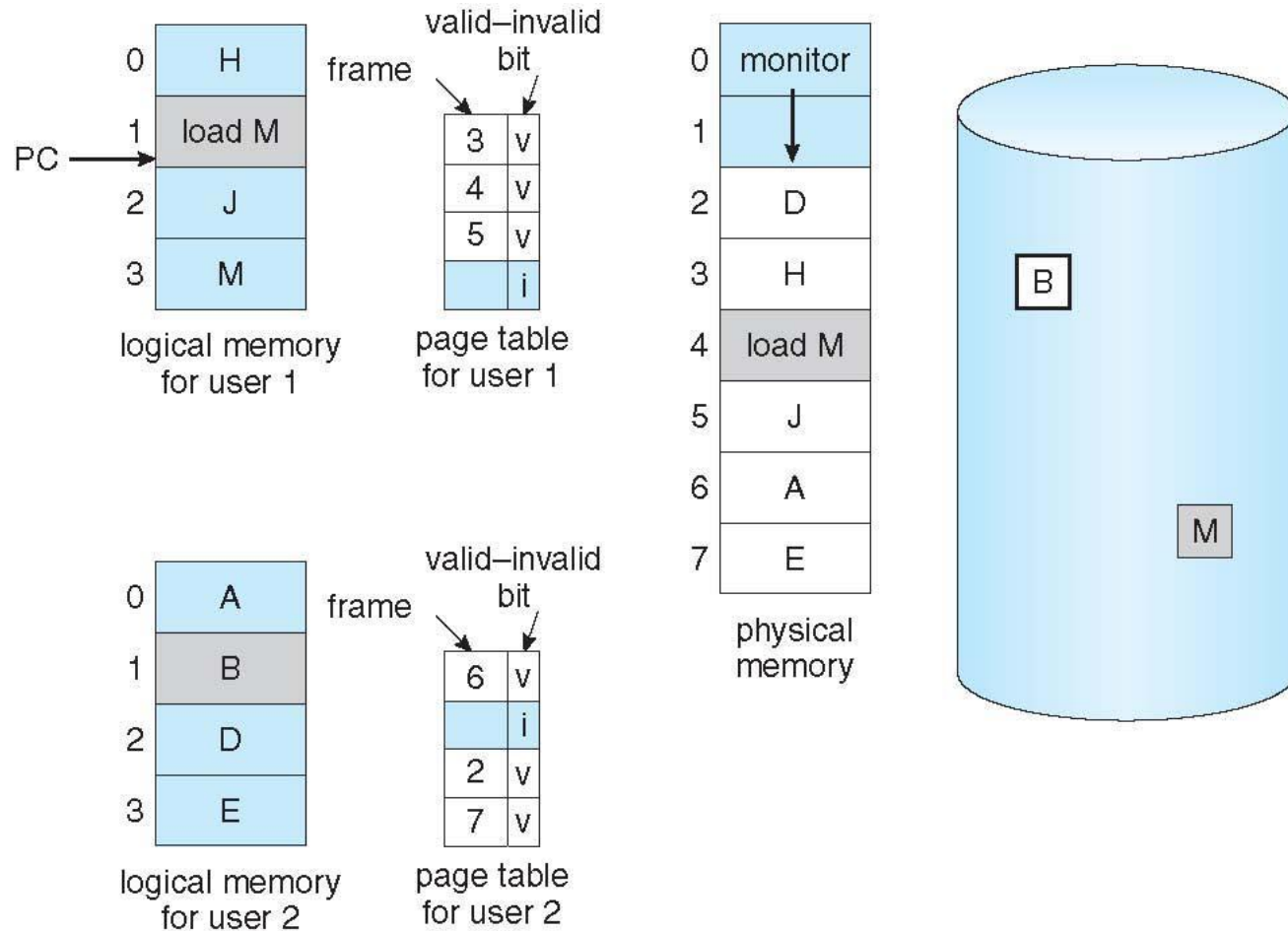- Put the process requesting more pages into a wait queue until some free frames become available.

# What Happens if There is no Free Frame?

- Swap some process out of memory completely, freeing up its page frames.

- Find some page in memory that isn't being used right now, and swap that page only out to disk, freeing up a frame that can be allocated to the process requesting it.

- This is known as page replacement, and is the most common solution.

- There are many different algorithms for page replacement, which is the subject of the remainder of this section.

# 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
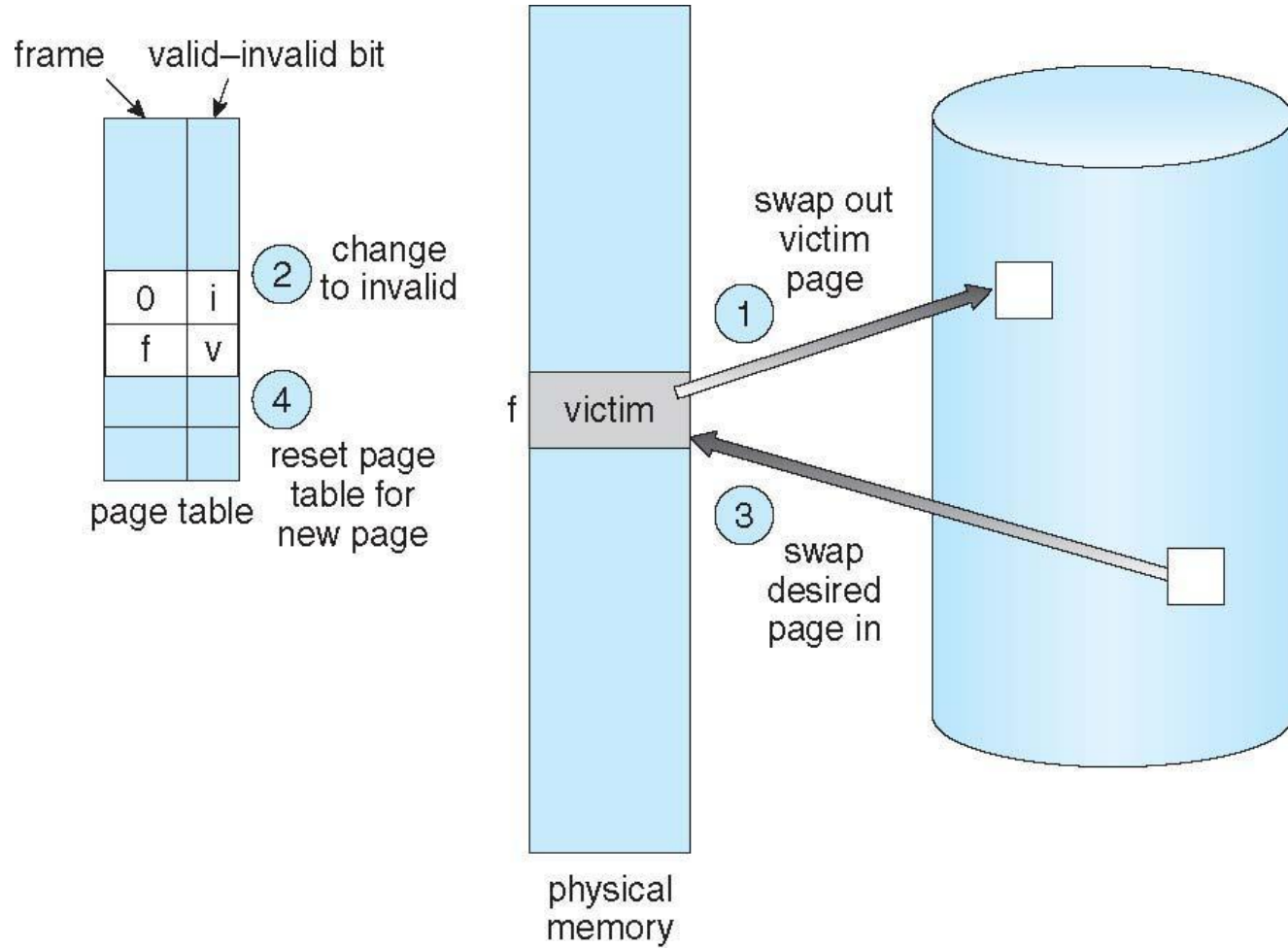
# Need For Page Replacement

# 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

# 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, Change all related page tables to indicate that this page is no longer in memory.

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



frame    valid–invalid bit

| 0 | i |
| f | v |

② change to invalid

④ reset page table for new page

page table

f  victim

physical memory

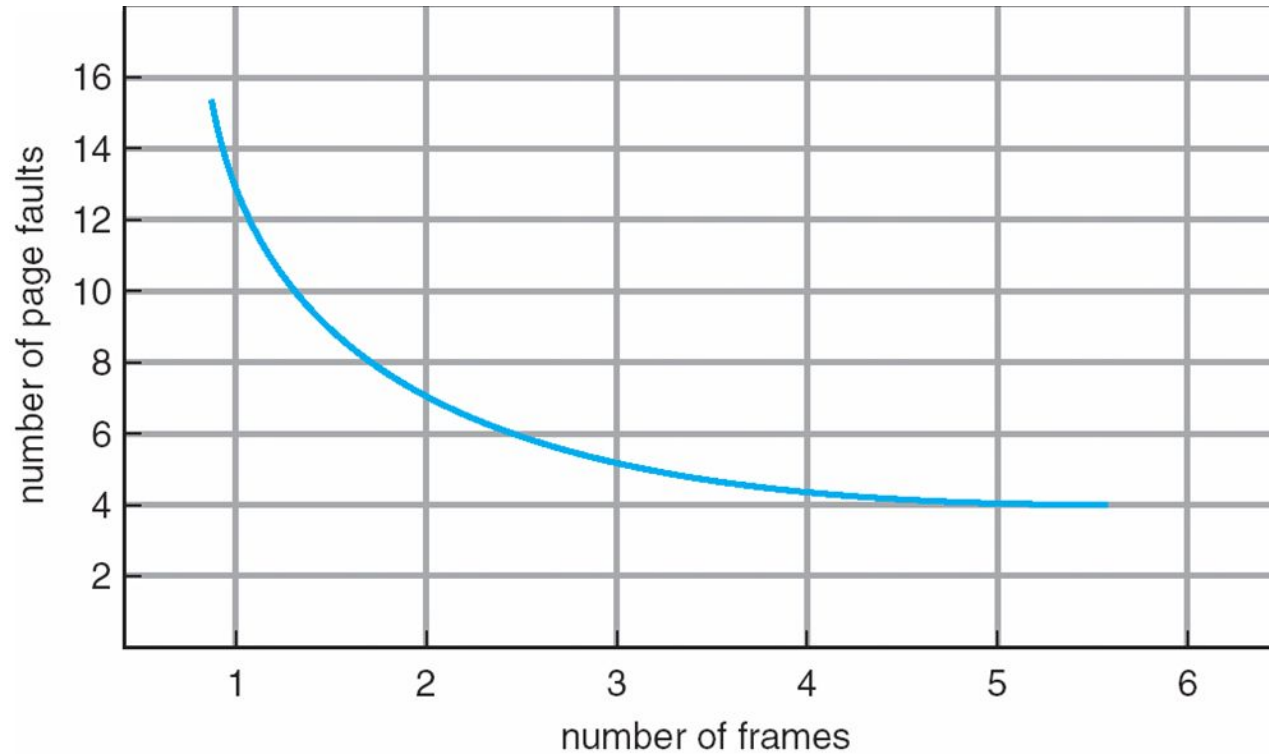① swap out victim page

③ swap desired page in

# 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
- The overall goal in selecting and tuning these algorithms is to generate the fewest number of overall page faults.
- Evaluate algorithm by running it on a particular string of memory references (reference string) and computing the number of page faults on that string
  - String is just page numbers, not full addresses
  - Repeated access to the same page does not cause a page fault
  - Results depend on number of frames available
- In all our examples, the **reference string** of referenced page numbers is

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

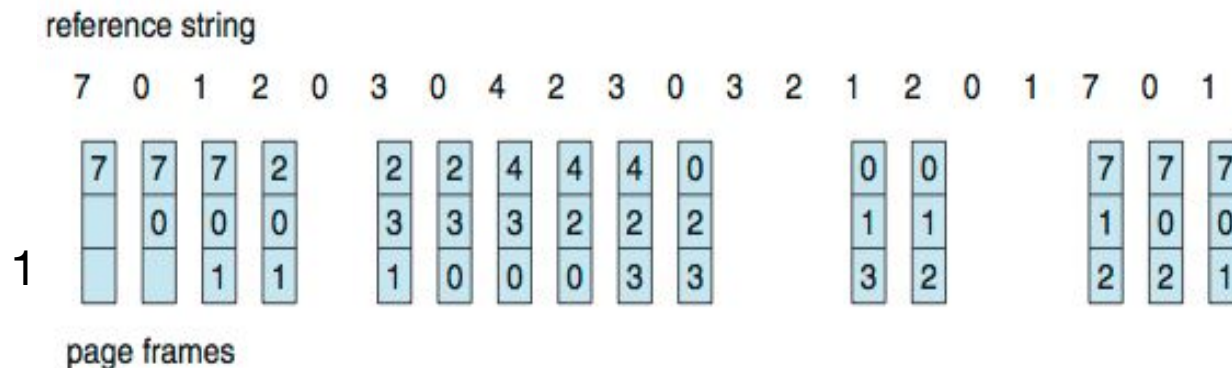# Graph of Page Faults Versus The Number of Frames

# FIFO Page Replacement

- A simple and obvious page replacement strategy is FIFO, i.e. first-in-first-out.

- As new pages are brought in, they are added to the tail of a queue, and the page at the head of the queue is the next victim.
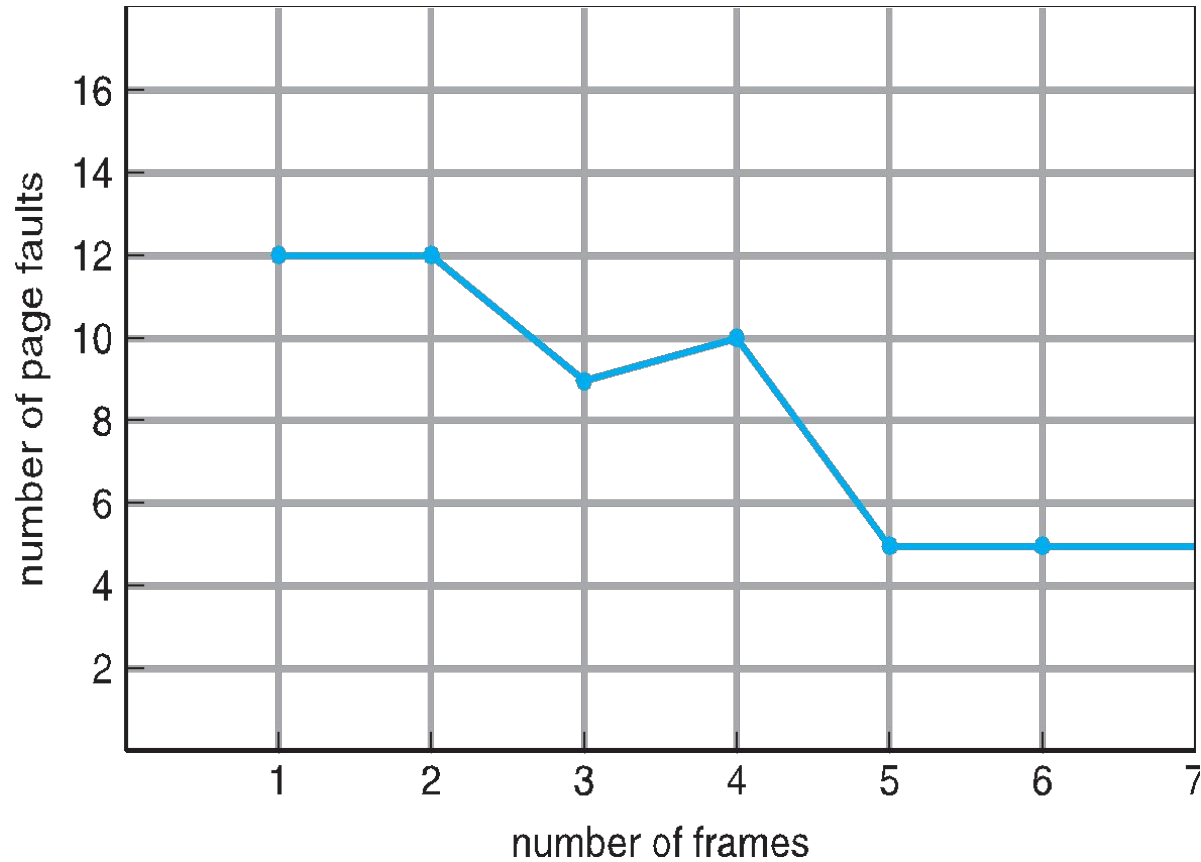
# First-In-First-Out (FIFO) Algorithm

- 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)



reference string

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

page frames

- 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 Illustrating Belady's Anomaly

# Optimal Page Replacement

- The discovery of Belady's anomaly lead to the search for an optimal page-replacement algorithm,
  - which is simply that which yields the lowest of all possible page-faults, and which does not suffer from Belady's anomaly.
- Such an algorithm does exist, and is called OPT or MIN.
- This algorithm is simply "Replace the page that will not be used for the longest time in the future."

# Optimal Algorithm

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

reference string

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



page frames

# Optimal Algorithm

- Unfortunately OPT cannot be implemented in practice, because it requires foretelling the future,

- In practice most page-replacement algorithms try to approximate OPT by predicting ( estimating ) in one fashion or another what page will not be used for the longest period of time.

- The basis of FIFO is the prediction that the page that was brought in the longest time ago is the one that will not be needed again for the longest future time,

- there are many other prediction methods, all striving to match the performance of OPT.

# LRU Page Replacement

- The prediction behind LRU, the Least Recently Used, algorithm is that the page that has not been used in the longest time is the one that will not be used again in the near future.

- Distinction between FIFO and LRU: The former looks at the oldest load time, and the latter looks at the oldest use time.

- Some view LRU as analogous to OPT, except looking backwards in time instead of forwards.

# 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 |



page frames

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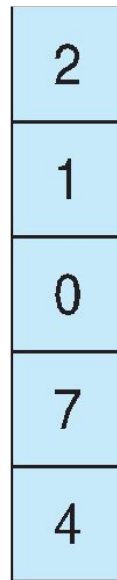- But how to implement?

# LRU Algorithm (Cont.)

- Counter implementation
  - Every page entry has a counter; every time page is referenced through this entry, copy the clock into the counter
  - When a page needs to be changed, look at the counters to find smallest value
    - Search through table needed

- Stack implementation
  - other approach is to use a stack,
  - whenever a page is accessed, pull that page from the middle of the stack and place it on the top.
  - The LRU page will always be at the bottom of the stack.
  - Because this requires removing objects from the middle of the stack, a doubly linked list is the recommended data structure.

- 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

reference string

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

| 2 |
|---|
| 1 |
| 0 |
| 7 |
| 4 |

stack
before
a

| 7 |
|---|
| 2 |
| 1 |
| 0 |
| 4 |

stack
after
b

a    b

# LRU Algorithm (Cont.)

- Note that both implementations of LRU require hardware support
  - either for incrementing the counter or for managing the stack
  - as these operations must be performed for every memory access.

- A stack algorithm is one in which the pages kept in memory for a frame set of size N will always be a subset of the pages kept for a frame size of N + 1.

- In the case of LRU, the top N pages of the stack will be the same for all frame set sizes of N or anything larger.

# Allocation of Frames

- Minimum Number of Frames
  - The absolute minimum number of frames that a process must be allocated is dependent on system architecture,
  - Corresponds to the worst-case scenario of the number of pages that could be touched by a single ( machine ) instruction.
  - If an instruction ( and its operands ) spans a page boundary, then multiple pages could be needed just for the instruction fetch.
  - Memory references in an instruction touch more pages, and if those memory locations can span page boundaries, then multiple pages could be needed for operand access also.

# Allocation of Frames

- Each process needs *minimum* number of frames
- Example:  IBM 370 – 6 pages to handle SS MOVE instruction:
  - instruction is 6 bytes, might span 2 pages
  - 2 pages to handle *from*
  - 2 pages to handle *to*
- *Maximum* of course is total frames in the system
- Two major allocation schemes
  - fixed allocation
  - priority allocation
- Many variations

# Fixed Allocation

- Equal allocation – For example, if there are 100 frames (after allocating frames for the OS) and 5 processes, give each process 20 frames
  - Keep some as free frame buffer pool

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

$s_i$ = size of process $p_i$

$S = \sum s_i$

$m$ = total number of frames

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

$m = 64$

$s_1 = 10$

$s_2 = 127$

$a_1 = \dfrac{10}{137} \times 62 \approx 4$

$a_2 = \dfrac{127}{137} \times 62 \approx 57$

# 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

# 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
- 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 **lgroups**
    - 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 lgroup

# Thrashing

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

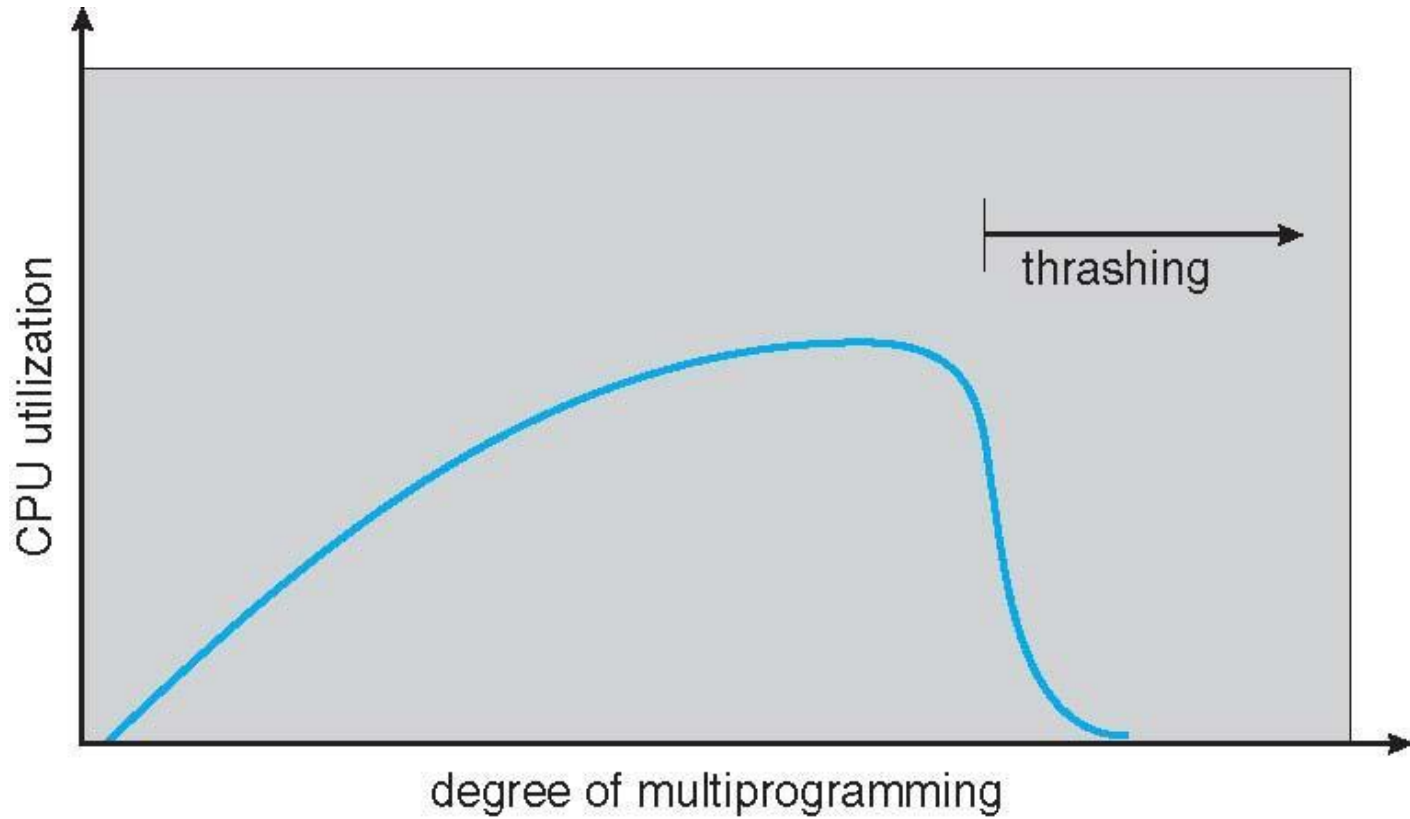- **Thrashing** ≡ a process is busy swapping pages in and out

# Cause of Thrashing

- Early process scheduling schemes would control the level of multiprogramming allowed based on CPU utilization, adding in more processes when CPU utilization was low.

- The problem is that when memory filled up and processes started spending lots of time waiting for their pages to page in, then CPU utilization would lower, causing the schedule to add in even more processes and exacerbating the problem!

- Eventually the system would essentially grind to a halt.

# Cause of Thrashing

- Local page replacement policies can prevent one thrashing process from taking pages away from other processes
  - but it still tends to clog up the I/O queue, thereby slowing down any other process that needs to do even a little bit of paging

# Thrashing (Cont.)

# Thrashing (Cont.)

- To prevent thrashing we must provide processes with as many frames as they really need "right now", but how do we know what that is?

- The locality model notes that processes typically access memory references in a given locality, making lots of references to the same general area of memory before moving periodically to a new locality

- If we could just keep as many frames as are involved in the current locality, then page faulting would occur primarily on switches from one locality to another

# Demand Paging and Thrashing

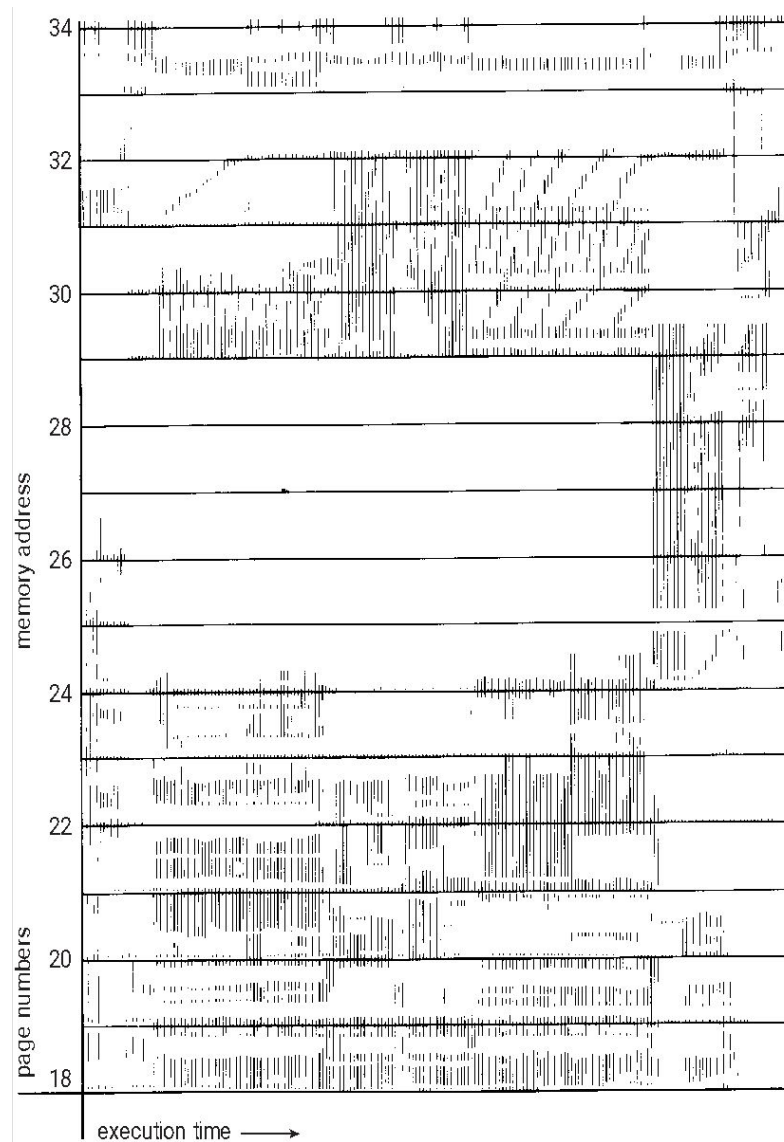- Why does demand paging work?
  **Locality model**
  - Process migrates from one locality to another
  - Localities may overlap

- Why does thrashing occur?
  $\Sigma$ size of locality > total memory size
  - Limit effects by using local or priority page replacement

# Locality In A Memory-Reference Pattern
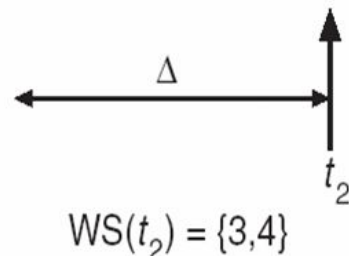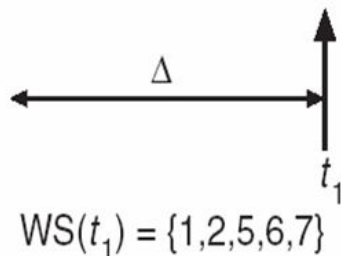
# Working-Set Model

- working set model is based on the concept of locality, and defines a working set window, of length delta.

- Whatever pages are included in the most recent delta page references are said to be in the processes working set window, and comprise its current working set,

# 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

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

page reference table

$$\ldots 2\ 6\ 1\ 5\ 7\ 7\ 7\ 7\ 5\ 1\ 6\ 2\ 3\ 4\ 1\ 2\ 3\ 4\ 4\ 4\ 3\ 4\ 3\ 4\ 4\ 4\ 1\ 3\ 2\ 3\ 4\ 4\ 4\ 3\ 4\ 4\ 4 \ldots$$

$\Delta$              $\Delta$

$t_1$                    $t_2$

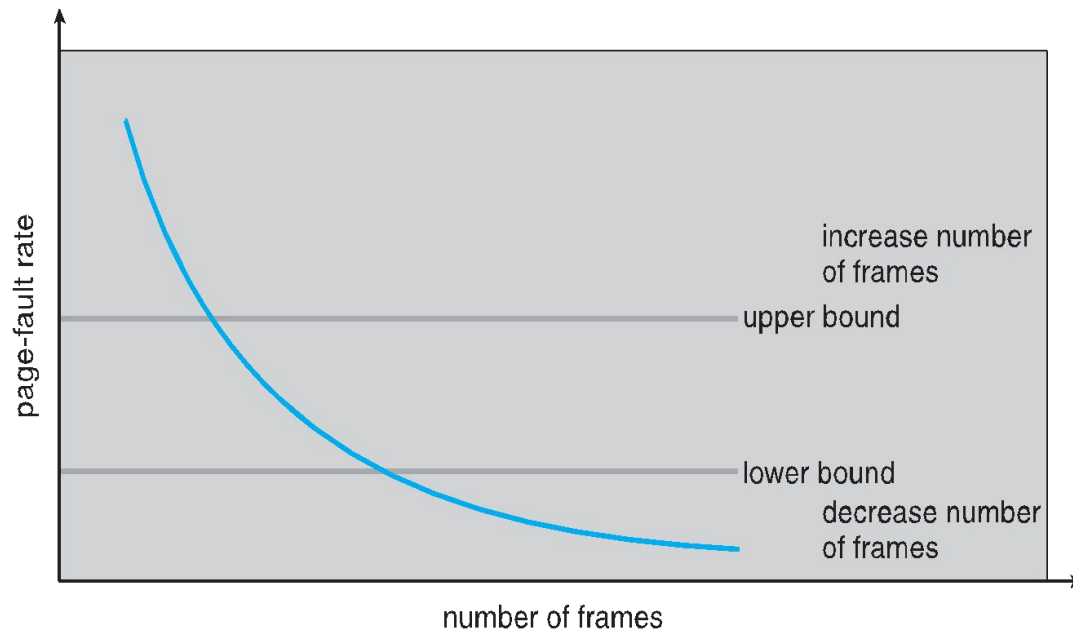$WS(t_1) = \{1,2,5,6,7\}$          $WS(t_2) = \{3,4\}$

# 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?
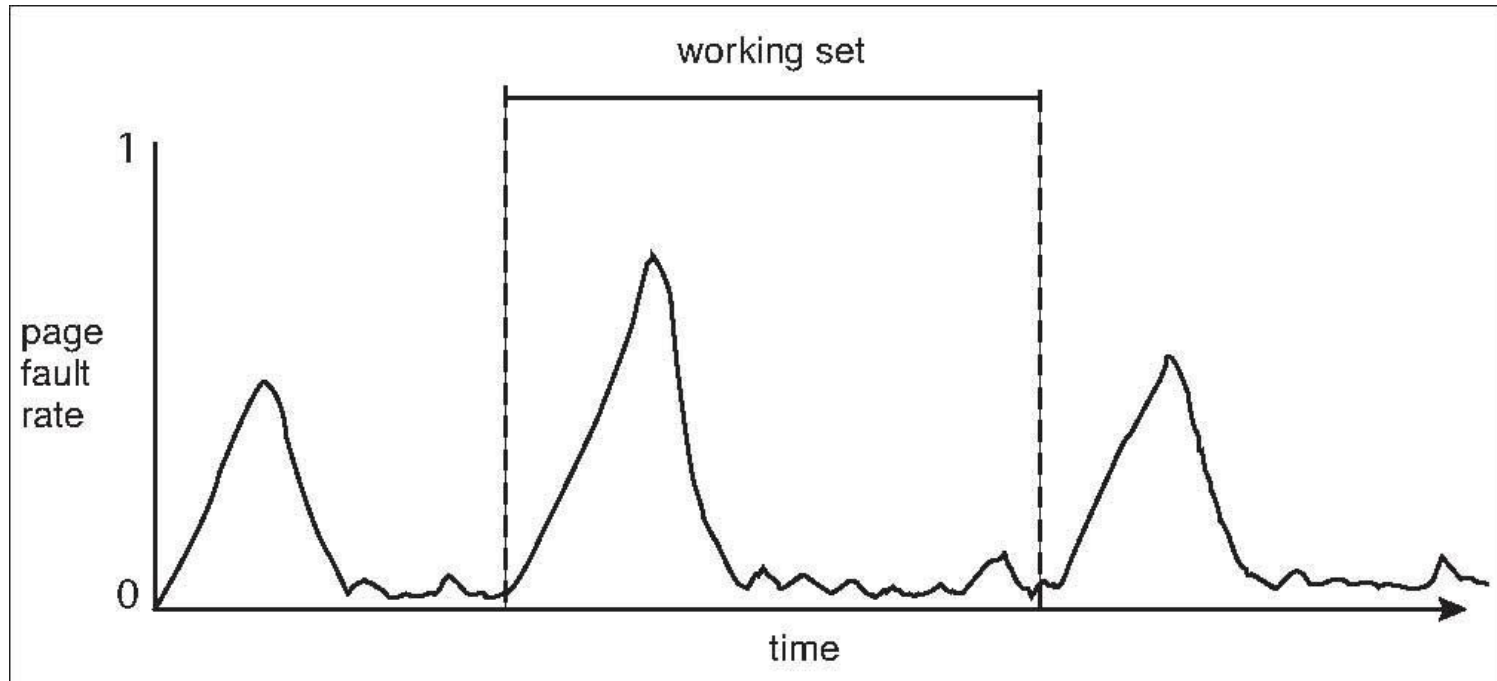- Improvement = 10 bits and interrupt every 1000 time units

# Page-Fault Frequency

- More direct approach than WSS
- Establish "acceptable" **page-fault frequency** (**PFF**) 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

- Direct relationship between working set of a process and its page-fault rate
- Working set changes over time
- Peaks and valleys over time

# End of Chapter 9