ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA
DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

# Real-Time Operating Systems M

## 7. Paging • Virtual Memory

# Notice

The course material includes slides downloaded from:

http://codex.cs.yale.edu/avi/os-book/

*(slides by Silberschatz, Galvin, and Gagne, associated with Operating System Concepts, 9th Edition, Wiley, 2013)*

and

http://retis.sssup.it/~giorgio/rts-MECS.html

*(slides by Buttazzo, associated with Hard Real-Time Computing Systems, 3rd Edition, Springer, 2011)*

which has been edited to suit the needs of this course.

The slides are authorized for personal use only.

Any other use, redistribution, and any for profit sale of the slides (in any form) requires the consent of the copyright owners.

# Implementation of Page Table

- Page table is kept in main memory

- **Page-table base register** (**PTBR**) points to the page table

- **Page-table length register** (**PTLR**) indicates size of the page table

- In this scheme every data/instruction access requires two memory accesses
    - One for the page table and one for the data / instruction

- The two memory access problem can be solved by the use of a special fast-lookup hardware cache called **associative memory** or **translation look-aside buffers** (**TLBs**)

- Some TLBs store **address-space identifiers** (**ASIDs**) in each TLB entry – uniquely identifies each process to provide address-space protection for that process
    - Otherwise need to flush at every context switch

- TLBs typically small (64 to 1,024 entries)

- On a TLB miss, value is loaded into the TLB for faster access next time
    - Replacement policies must be considered
    - Some entries can be **wired down** for permanent fast access

# Associative Memory

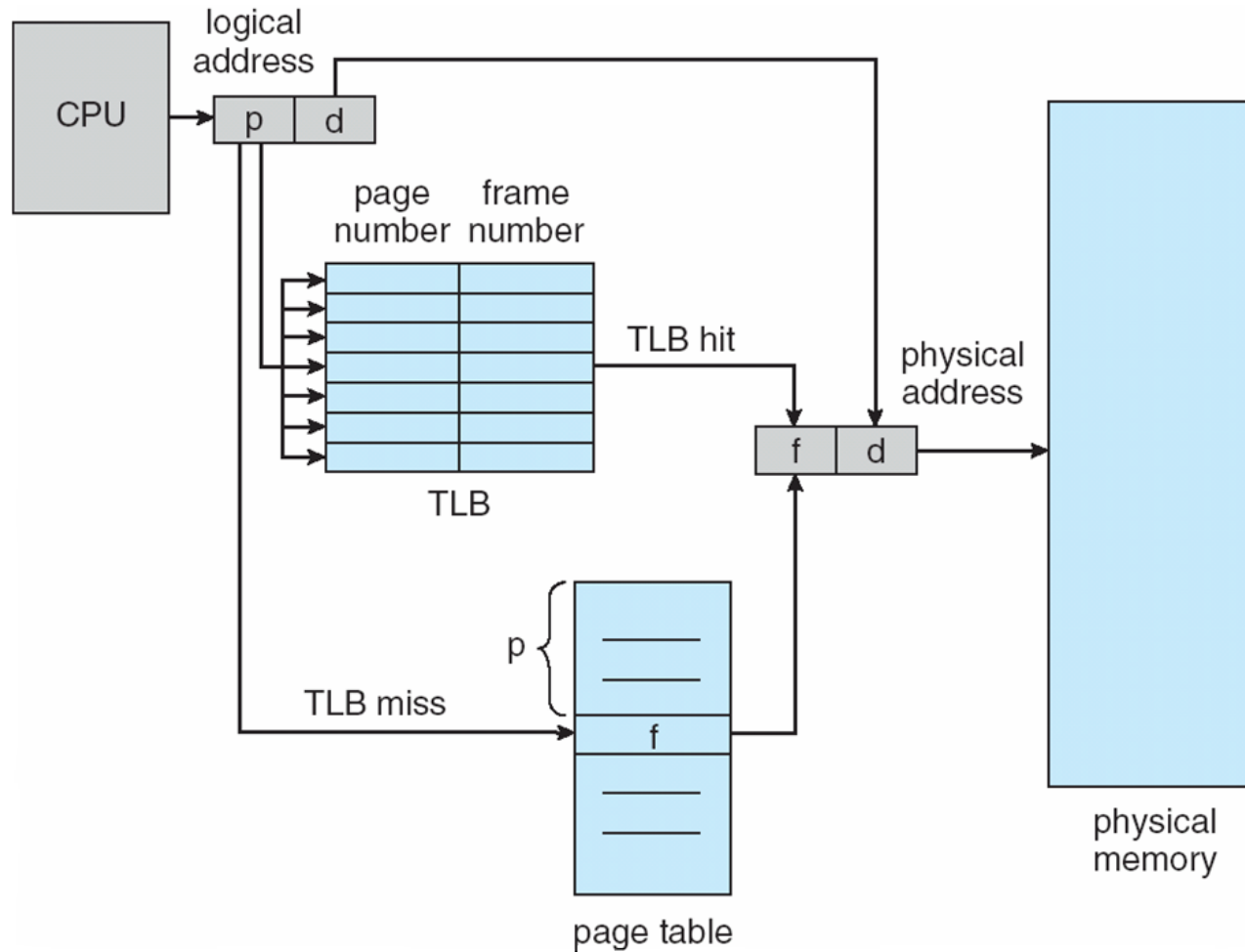- Associative memory – parallel search

| Page # | Frame # |
|--------|---------|
|        |         |
|        |         |
|        |         |
|        |         |

- Address translation (p, d)
  - If p is in associative register, get frame # out
  - Otherwise get frame # from page table in memory

# Paging Hardware With TLB

# Effective Access Time

- Associative Lookup = $\varepsilon$ time unit
  - Can be < 10% of memory access time

- Hit ratio = $\alpha$
  - Hit ratio – percentage of times that a page number is found in the associative registers; ratio related to number of associative registers

- Consider $\alpha$ = 80%, $\varepsilon$ = 20ns for TLB search, 100ns for memory access

- **Effective Access Time** (**EAT**)

$$EAT = (1 + \varepsilon)\, \alpha + (2 + \varepsilon)(1 - \alpha)$$
$$= 2 + \varepsilon - \alpha$$

- Consider $\alpha$ = 80%, $\varepsilon$ = 20ns for TLB search, 100ns for memory access
  - EAT = 0.80 x 120 + 0.20 x 220 = 140ns

- Consider more realistic hit ratio -> $\alpha$ = 99%, $\varepsilon$ = 20ns for TLB search, 100ns for memory access
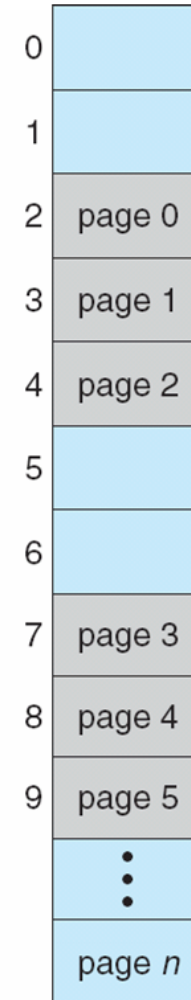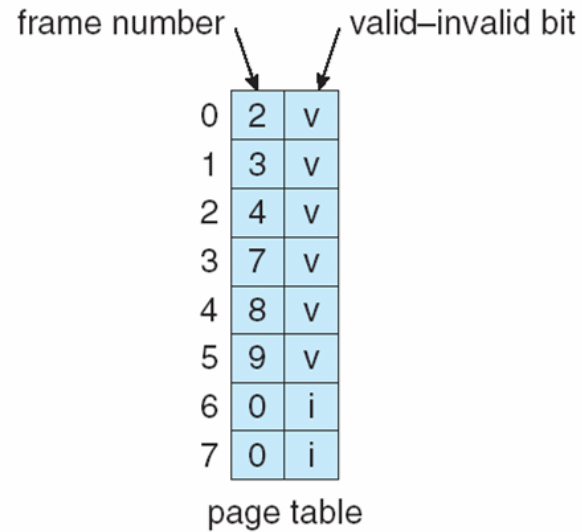  - EAT = 0.99 x 120 + 0.01 x 220 = 121ns

# Memory Protection
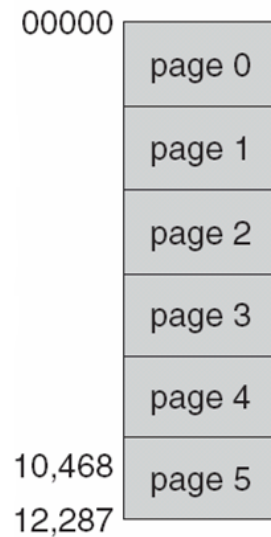
- Memory protection implemented by associating protection bit with each frame to indicate if read-only or read-write access is allowed
  - Can also add more bits to indicate page execute-only, and so on

- **Valid-invalid** bit attached to each entry in the page table:
  - "valid" indicates that the associated page is in the process' logical address space, and is thus a legal page
  - "invalid" indicates that the page is not in the process' logical address space
  - Or use **page-table length register** (**PTLR**)

- Any violations result in a trap to the kernel

# Valid (v) or Invalid (i) Bit In A Page Table

```
00000
        ┌──────────┐
        │  page 0  │
        ├──────────┤
        │  page 1  │
        ├──────────┤
        │  page 2  │
        ├──────────┤
        │  page 3  │
        ├──────────┤
        │  page 4  │
10,468  ├──────────┤
        │  page 5  │
12,287  └──────────┘
```

frame number          valid–invalid bit

```
    ┌───┬───┐
0   │ 2 │ v │
1   │ 3 │ v │
2   │ 4 │ v │
3   │ 7 │ v │
4   │ 8 │ v │
5   │ 9 │ v │
6   │ 0 │ i │
7   │ 0 │ i │
    └───┴───┘
  page table
```

```
0  ┌──────────┐
   │          │
1  ├──────────┤
   │          │
2  ├──────────┤
   │  page 0  │
3  ├──────────┤
   │  page 1  │
4  ├──────────┤
   │  page 2  │
5  ├──────────┤
   │          │
6  ├──────────┤
   │          │
7  ├──────────┤
   │  page 3  │
8  ├──────────┤
   │  page 4  │
9  ├──────────┤
   │  page 5  │
   ├──────────┤
   │    •     │
   │    •     │
   │    •     │
   ├──────────┤
   │  page n  │
   └──────────┘
```

# Shared Pages

- **Shared code**

  - One copy of read-only (**reentrant**) code shared among processes (i.e., text editors, compilers, window systems)

  - Similar to multiple threads sharing the same process space

  - Also useful for interprocess communication if sharing of read-write pages is allowed
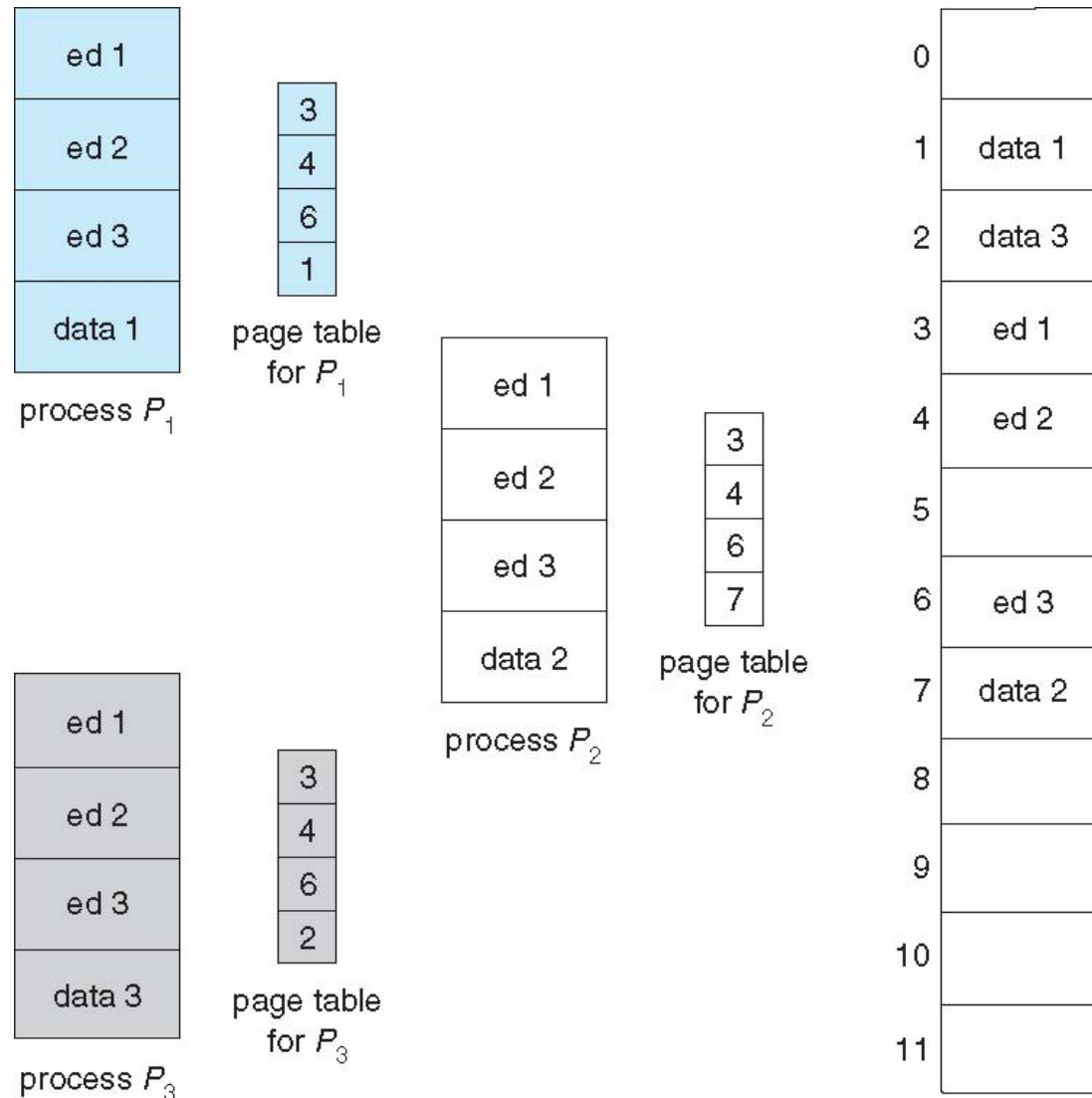
- **Private code and data**

  - Each process keeps a separate copy of the code and data

  - The pages for the private code and data can appear anywhere in the logical address space

# Shared Pages Example

# Structure of the Page Table

- Memory structures for paging can get huge using straight-forward methods
  - Consider a 32-bit logical address space as on modern computers
  - Page size of 4 KB ($2^{12}$)
  - Page table would have 1 million entries ($2^{32} / 2^{12}$)
  - If each entry is 4 bytes -> 4 MB of physical address space / memory for page table alone
    - That amount of memory used to cost a lot
    - Don't want to allocate that contiguously in main memory

- Hierarchical Paging

- Hashed Page Tables

- Inverted Page Tables

# Hierarchical Page Tables

- Break up the logical address space into multiple page tables

- A simple technique is a two-level page table

- We then page the page table

# Two-Level Page-Table Scheme

# Two-Level Paging Example

- A logical address (on 32-bit machine with 4K page size) is divided into:
  - a page number consisting of 20 bits
  - a page offset consisting of 12 bits

- Since the page table is paged, the page number is further divided into:
  - a 10-bit page number
  - a 10-bit page offset

- Thus, a logical address is as follows:

| page number | | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 10 | 10 | 12 |

- where $p_1$ is an index into the outer page table, and $p_2$ is the displacement within the page of the inner page table
- Known as **forward-mapped page table**

# Address-Translation Scheme



logical address

| $p_1$ | $p_2$ | d |

$p_1$ → outer page table

$p_2$ → page of page table

d →

# 64-bit Logical Address Space

- Even two-level paging scheme not sufficient
- If page size is 4 KB ($2^{12}$)
  - Then page table has $2^{52}$ entries
  - If two level scheme, inner page tables could be $2^{10}$ 4-byte entries
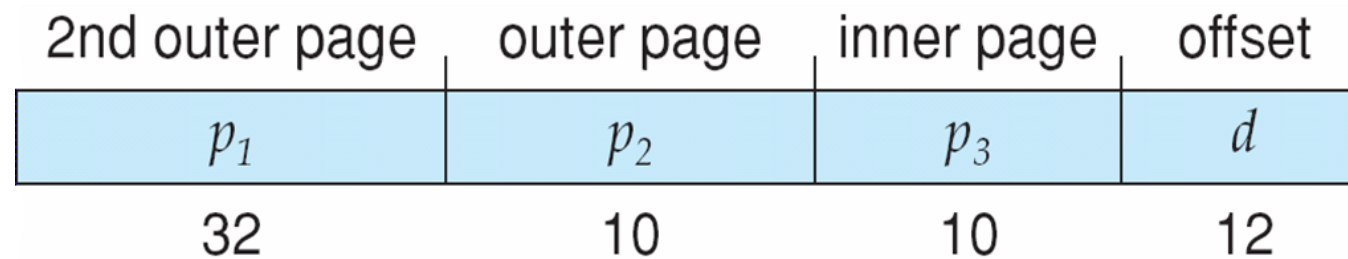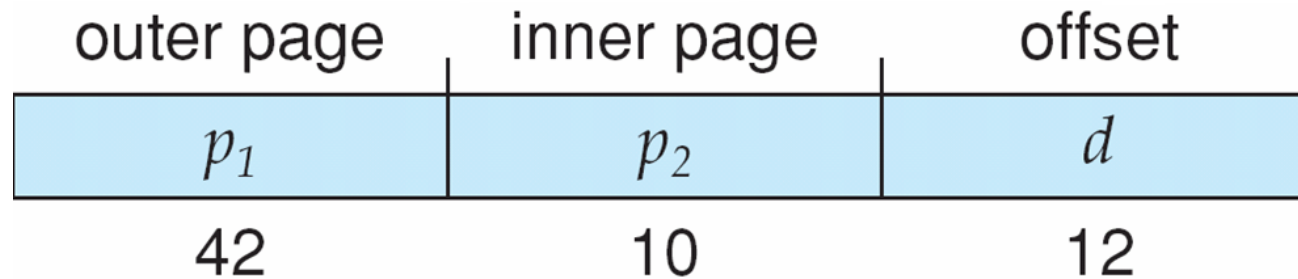  - Address would look like

| outer page | inner page | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

  - Outer page table has $2^{42}$ entries or $2^{44}$ bytes
  - One solution is to add a 2$^{nd}$ outer page table
  - But in the following example the 2$^{nd}$ outer page table is still $2^{34}$ bytes in size
    - And possibly 4 memory access to get to one physical memory location

# Three-level Paging Scheme

| outer page | inner page | offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 42 | 10 | 12 |

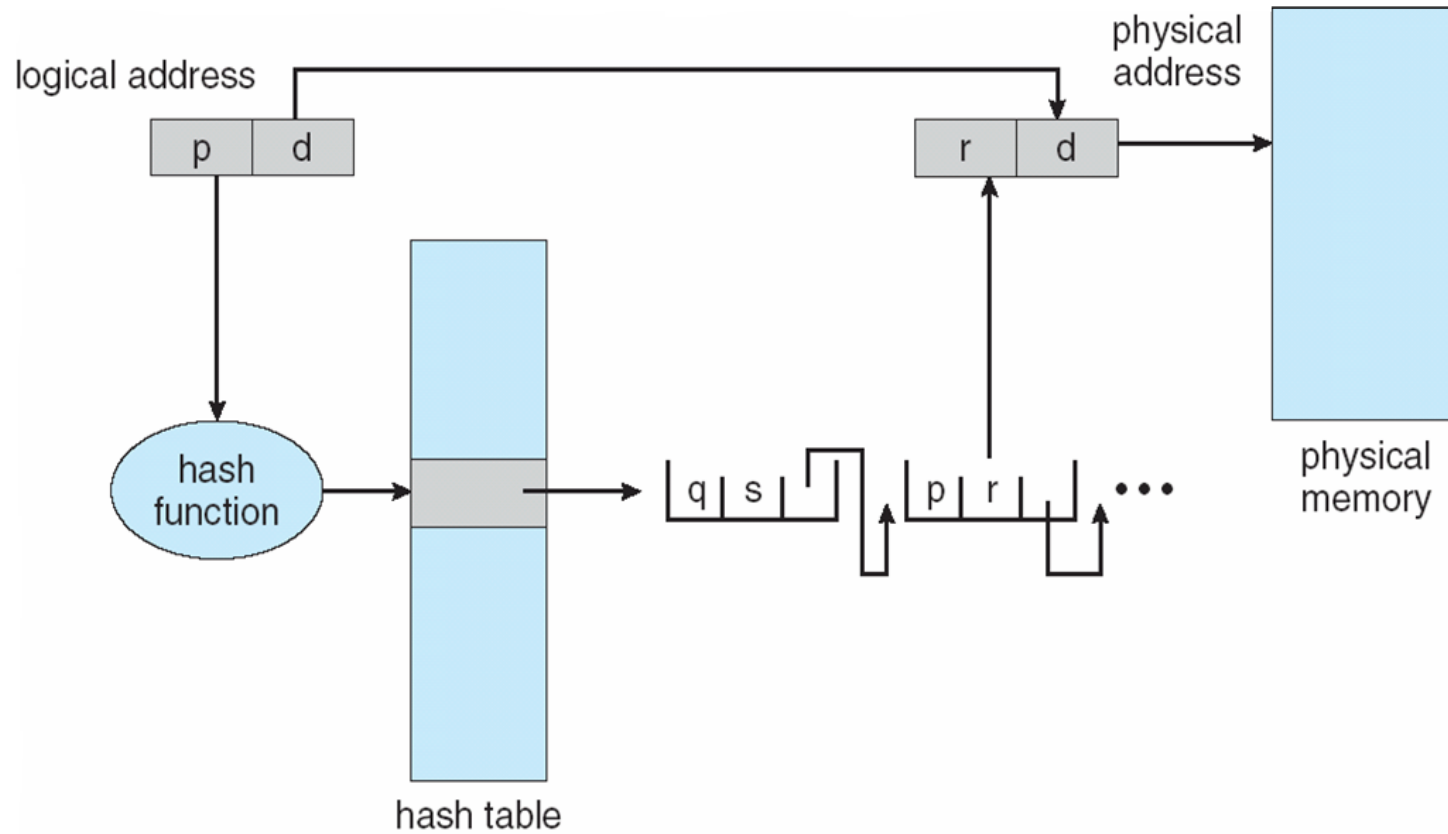| 2nd outer page | outer page | inner page | offset |
|:---:|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $p_3$ | $d$ |
| 32 | 10 | 10 | 12 |

# Hashed Page Tables

- Common in address spaces > 32 bits

- The virtual page number is hashed into a page table
  - This page table contains a chain of elements hashing to the same location

- Each element contains (1) the virtual page number (2) the value of the mapped page frame (3) a pointer to the next element

- Virtual page numbers are compared in this chain searching for a match
  - If a match is found, the corresponding physical frame is extracted

- Variation for 64-bit addresses is **clustered page tables**
  - Similar to hashed but each entry refers to several pages (such as 16) rather than 1
  - Especially useful for **sparse** address spaces (where memory references are non-contiguous and scattered)
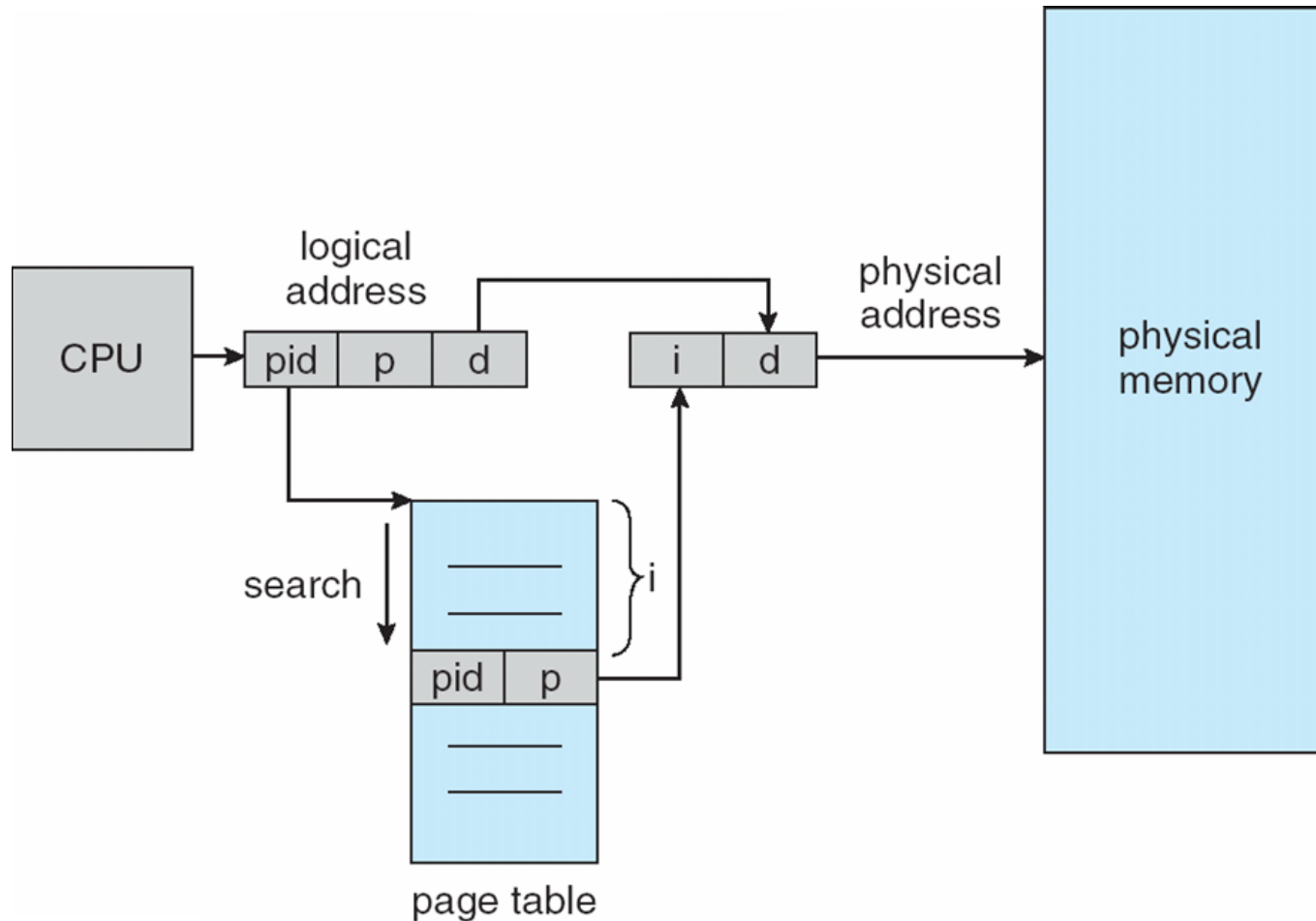
# Hashed Page Table

# Inverted Page Table

■ Rather than each process having a page table and keeping track of all possible logical pages, track all physical pages

■ One entry for each real page of memory

■ Entry consists of the virtual address of the page stored in that real memory location, with information about the process that owns that page

■ Decreases memory needed to store each page table, but increases time needed to search the table when a page reference occurs

■ Use hash table to limit the search to one — or at most a few — page-table entries
  ● TLB can accelerate access

■ But how to implement shared memory?
  ● One mapping of a virtual address to the shared physical address

# Inverted Page Table Architecture

# Example: Intel 32 Architecture

- Dominant industry chips

- Pentium CPUs are 32-bit and called IA-32 architecture

- Current Intel CPUs are 64-bit and called IA-64 architecture

- Many variations in the chips, main ideas in the textbook

# Example: The Intel IA-32 Architecture

- Supports both segmentation and segmentation with paging
  - Each segment can be 4 GB
  - Up to 16 K segments per process
  - Divided into two partitions
    - First partition of up to 8 K segments are private to process (kept in **local descriptor table** (**LDT**))
    - Second partition of up to 8K segments shared among all processes (kept in **global descriptor table** (**GDT**))

- CPU generates logical address
  - Selector given to segmentation unit
    - Which produces linear addresses

| $s$ | $g$ | $p$ |
|-----|-----|-----|
| 13  | 1   | 2   |

  - it
    - Which generates physical address in main memory
    - Paging units form equivalent of MMU
    - Pages sizes can be 4 KB or 4 MB

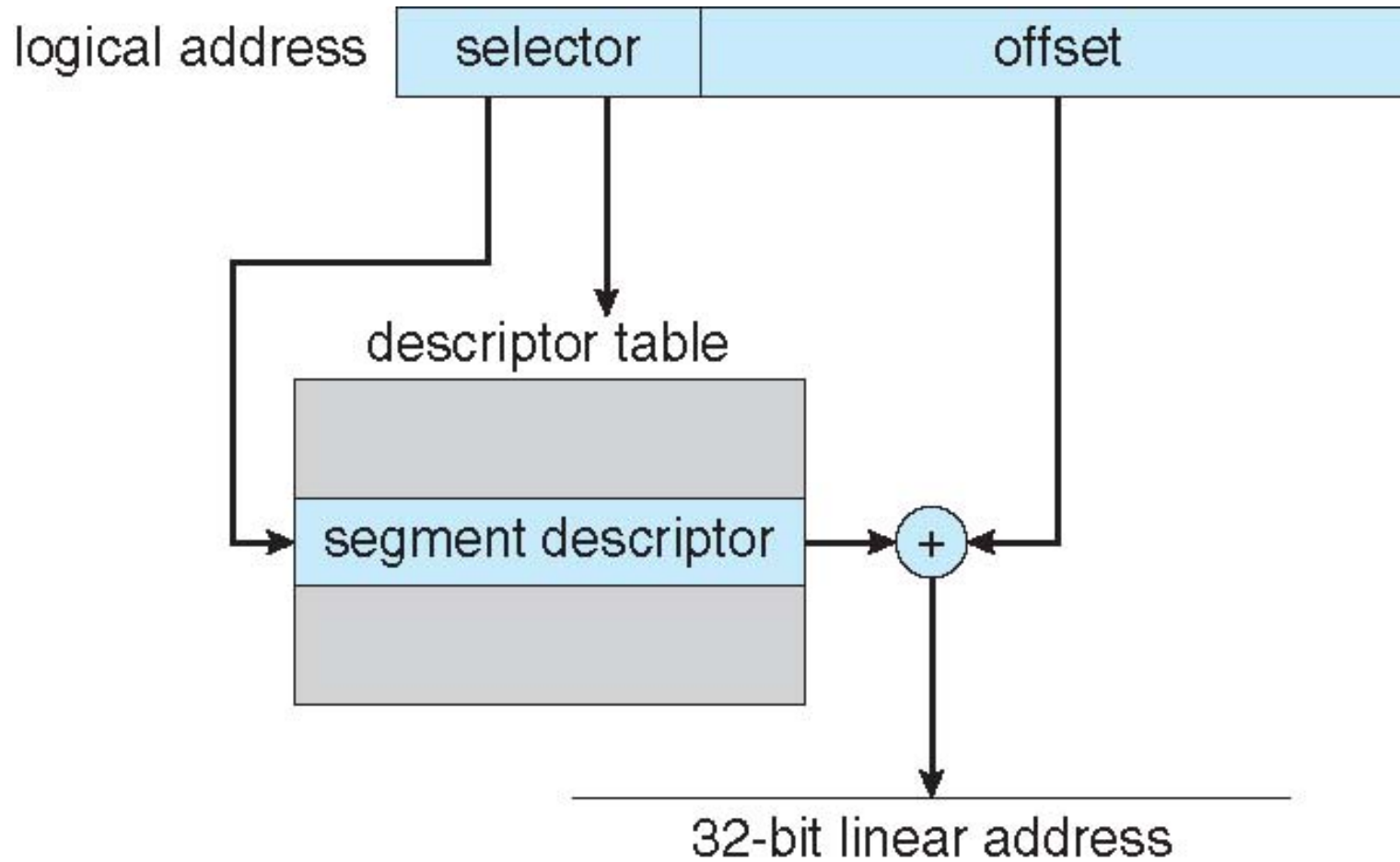**Silberschatz, Galvin and Gagne ©2013**

# Logical to Physical Address Translation in IA-32



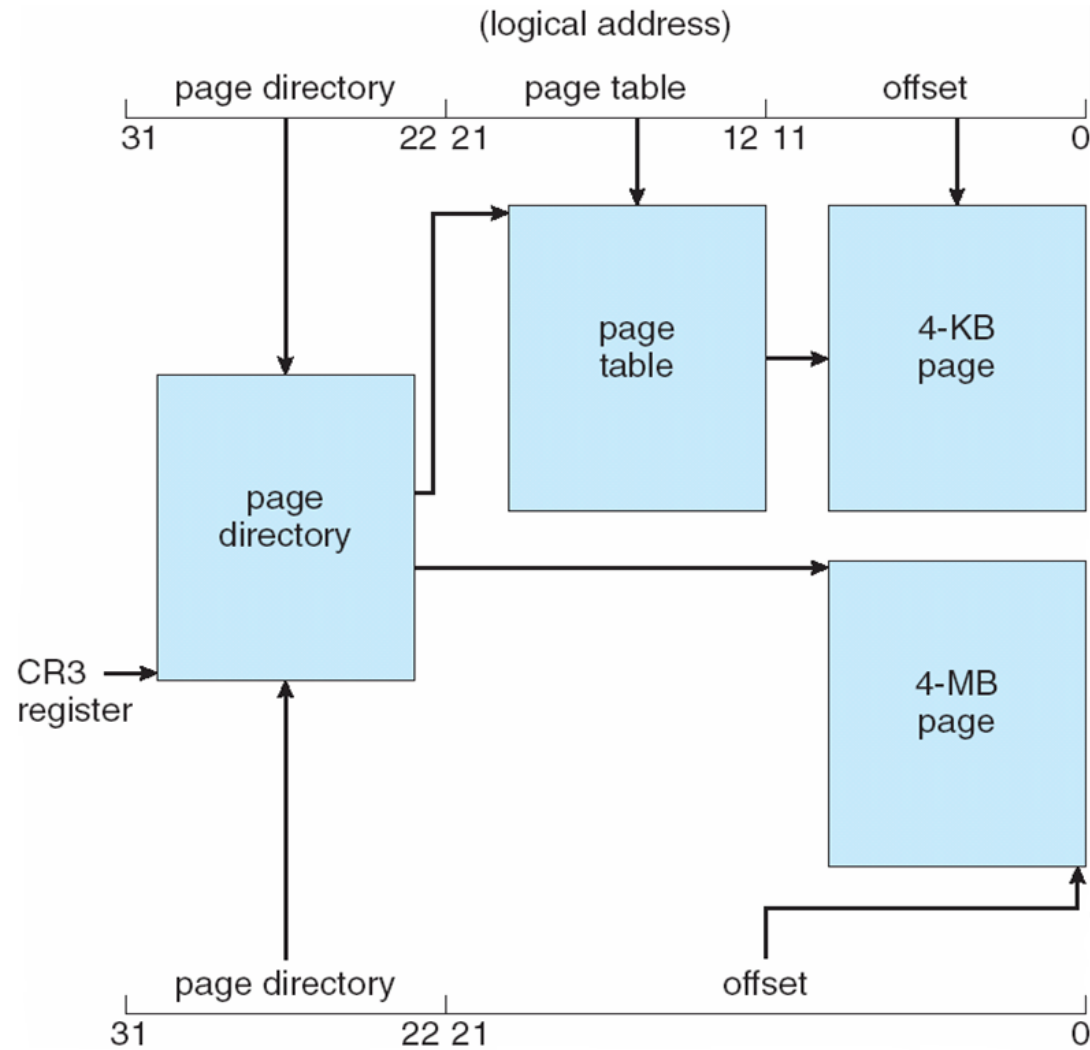| page number | | page offset |
|:---:|:---:|:---:|
| $p_1$ | $p_2$ | $d$ |
| 10 | 10 | 12 |

# Intel IA-32 Segmentation

# Intel IA-32 Paging Architecture

# Quizzes

- **Binding** of instructions and data to memory addresses is done at **load time**
- Segmentation and paging are two **alternative** ways to manage the main memory
- Instructions to be **executed** need to be in the main memory
- The **size of a program** is limited by the size of the physical memory
- A process **cannot be executed** if its logical address space is bigger than the physical address space
- An **inverted page table** contains, in each entry, a page number and a frame number
- A **hashed page table** contains, in each entry, a page number and a frame number
- **Paging** permits pages to be of arbitrary size
- **Segmentation** permits segments to be of arbitrary size
- A **segment table length register** (STLR) can be used to protect memory from illegal accesses
- A **page table length register** (PTLR) can be used to protect memory from illegal accesses
- A **page table base register** (PTBR) can be used to protect memory from illegal accesses
- **Swapping** moves out of the main memory **the entire image** of a process, in order to make space for the image of another process
- External fragmentation **implies** internal fragmentation
- Internal fragmentation **implies** external fragmentation
- An advantage of paging is the possibility of **sharing** common **code** or **data**.

# Exercise

- Consider a paging system with the page table stored in memory.

  - If a memory reference takes 200ns, how long does a paged memory refrence take?

  - If we add TLBs, and 75% of all page-table references are found in the TLBs, what is the effective memory reference time? (assume 0 time to find a TLB page entry, in case of TLB hit)

# Exercise

- Consider a logical address space of 64 pages of 1K words each, mapped onto a physical memory of 32 frames.

  - How many bits are there in the logical address?

  - How many bits are there in the physical address?

# Chapter 9: Virtual Memory

- Background
- Demand Paging
- Copy-on-Write
- Page Replacement
- Allocation of Frames
- Thrashing
- Memory-Mapped Files
- Allocating Kernel Memory
- Other Considerations
- Operating-System Examples

# 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

# 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

- **Virtual memory** – separation of user logical memory from physical memory
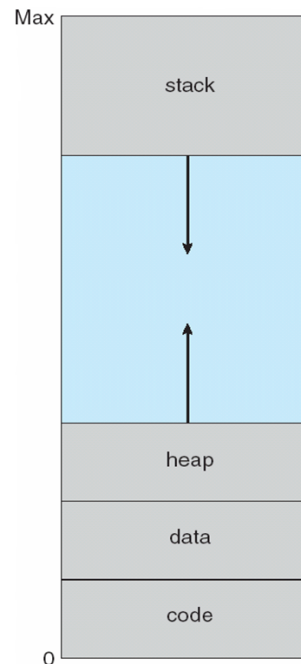
# Virtual Memory That is
# Larger Than Physical Memory

- Only part of the program needs to be in memory for execution

page 0
page 1
page 2

⋮

page v

virtual
memory

memory
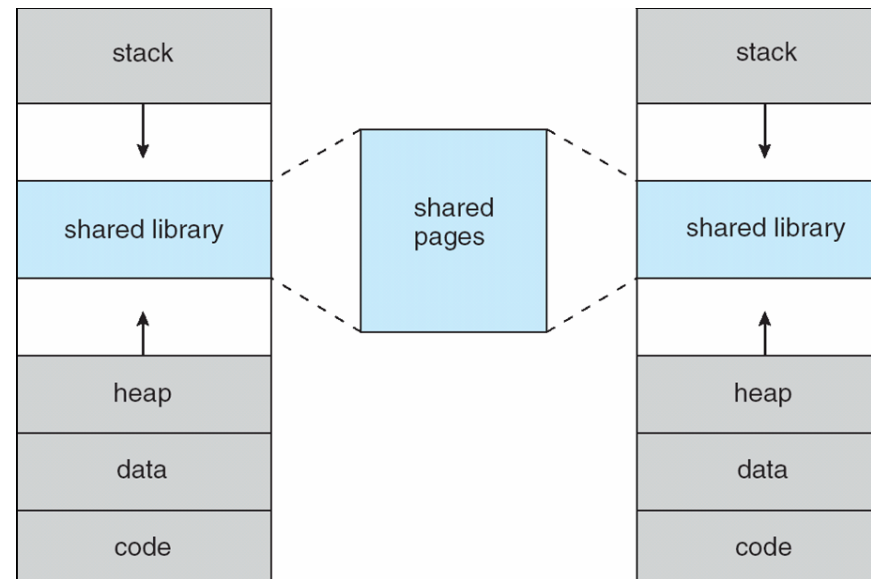map

physical
memory

# Virtual Address Space

- Logical address space can be much larger than physical address space
  - Easier to program
  - Enables **sparse** address spaces with holes left for growth, dynamically linked libraries, etc

# Virtual Address Space

- 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

# Virtual Memory

- **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
        - Easier to program
    - Allows address spaces to be shared by several processes
    - Allows for more efficient process creation
    - More programs running concurrently
        - More CPU utilization
    - Less I/O needed to load or swap processes

- Implementation of virtual memory by **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

- 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**
  - ["Swapper" moves entire process image]

# 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:

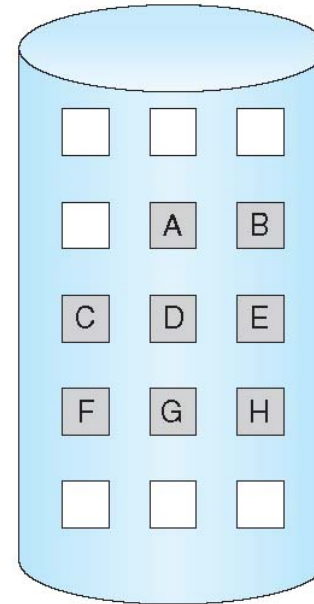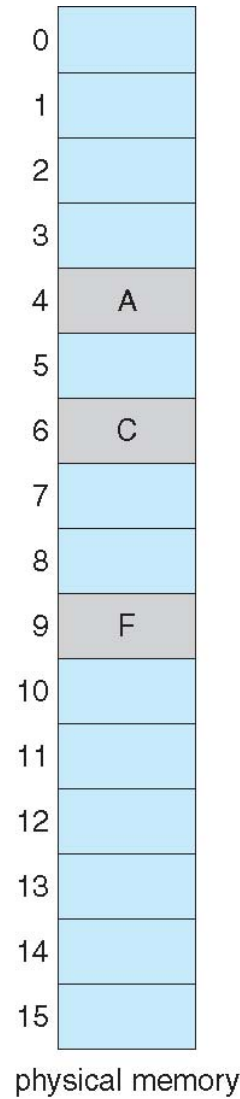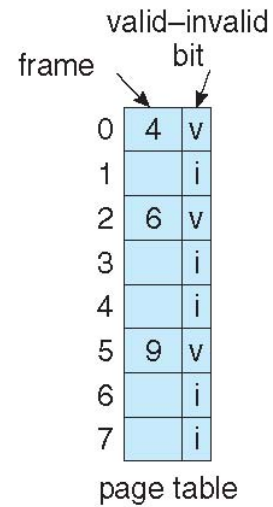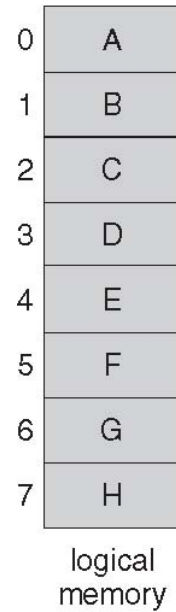| Frame # | valid-invalid bit |
|---------|:---:|
| | **v** |
| | **v** |
| | **v** |
| | **v** |
| | **i** |
| …. | |
| | **i** |
| | **i** |

page table

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

# Page Fault

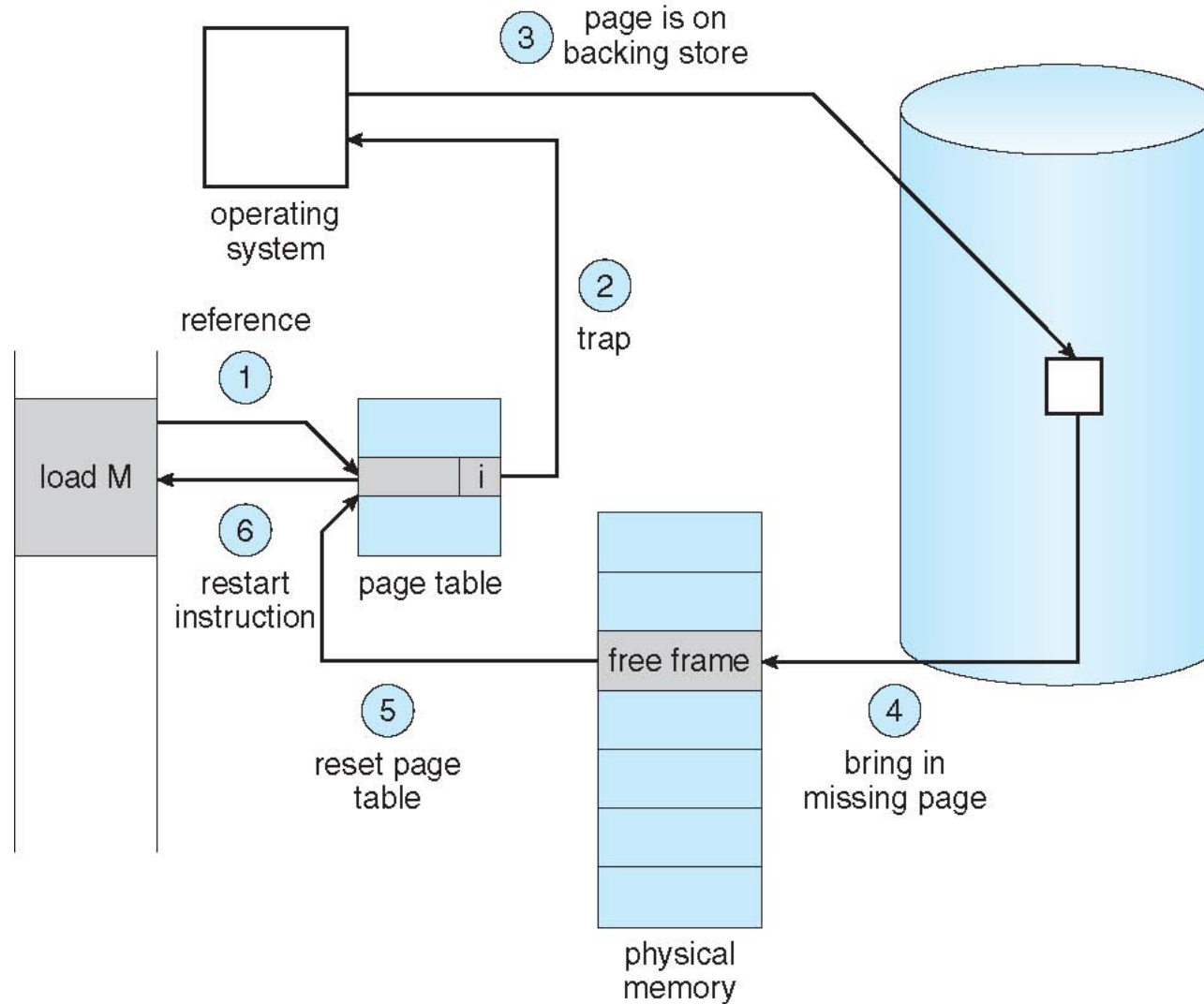- If there is a reference to a page that is not in memory, first reference to that page will trap to operating system:

    **page fault**

1. Operating system looks at internal table (in PCB) to check if legal reference:

    - Invalid reference ⇒ abort

    - Just not in memory ⇒ page in…

2. Get empty frame

3. Swap page into frame via scheduled disk operation

4. Reset tables to indicate page now in memory
   Set validation bit = **v**

5. Restart the instruction that caused the page fault

# Steps in Handling a Page Fault
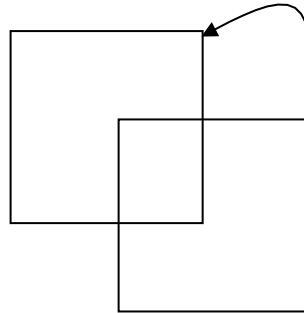
# Aspects of Demand Paging

- 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**
- 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



- Restart the whole operation?
  - What if source and destination overlap?
  - What if block straddles page boundary?
    - ▸ either use buffer to keep old values
    - ▸ or force page fault at beginning of operation
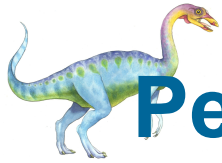    - ▸ Architecture/microcode must support proper restart (by being aware of page faults)

# Performance of Demand Paging

■ Stages in Demand Paging

1. Trap to the operating system
2. Save the user registers and process state
3. Determine that the interrupt was a page fault
4. Check that the page reference was legal and determine the location of the page on the disk
5. 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 (if CPU scheduling)
7. Receive an interrupt from the disk I/O subsystem (I/O completed)
8. Save the registers and process state for the other user
9. Determine that the interrupt was from the disk
10. Correct the page table and other tables to show page is now in memory
11. Wait for the CPU to be allocated to this process again
12. Restore the user registers, process state, and new page table, and then resume the interrupted instruction

# Performance of Demand Paging (Cont.)

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

- Effective Access Time (EAT)

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

# Demand Paging Example

- Memory access time = 200 nanoseconds

- Average page-fault service time = 8 milliseconds

- EAT = $(1 - p) \times 200 + p$ (8 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 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

# 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

- Prepage at startup all/some pages a process will need, before they are referenced
    - Reduces large number of page faults that occurs at process startup
    - But if prepaged pages are unused, I/O and memory was wasted
    - Assume $s$ pages are prepaged and $a$ of the pages is used
        - Is cost of $s * a$ save pages faults > or < than the cost of prepaging $s * (1- a)$ unnecessary pages?
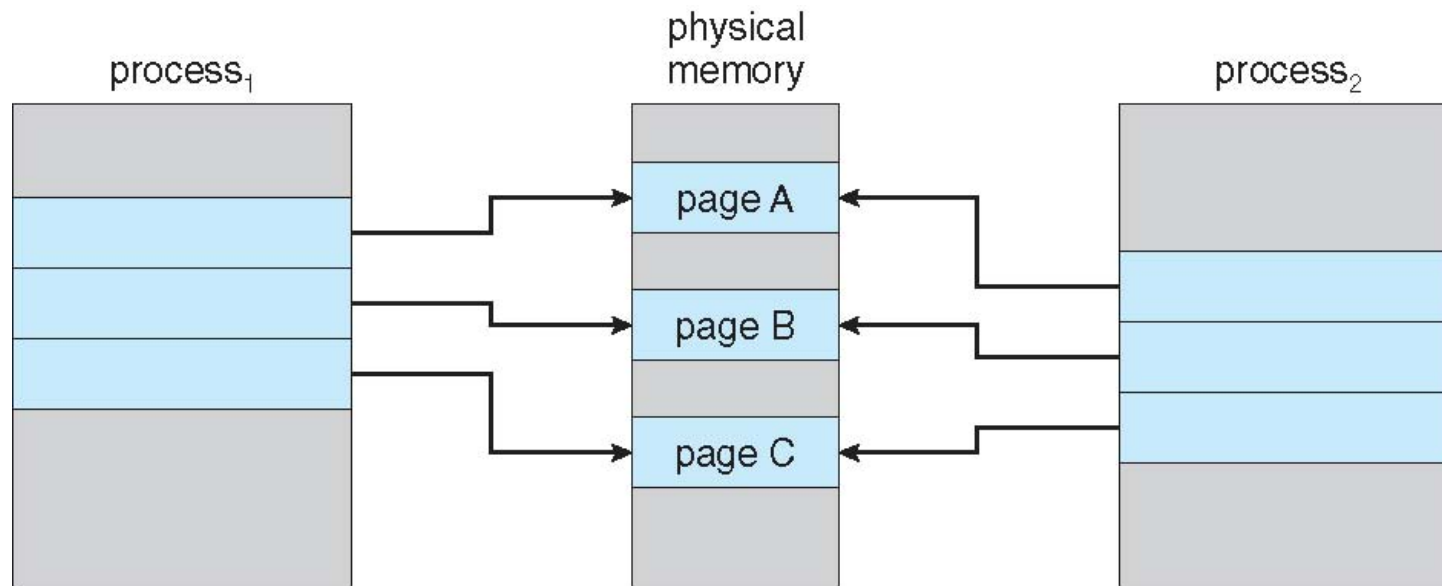        - $a$ near zero $\Rightarrow$ prepaging loses

# 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
    - Why zero-out a page before allocating it?
- `vfork()` variation on `fork()` system call has parent suspend and child using address space of parent
    - Designed to have child call `exec()`
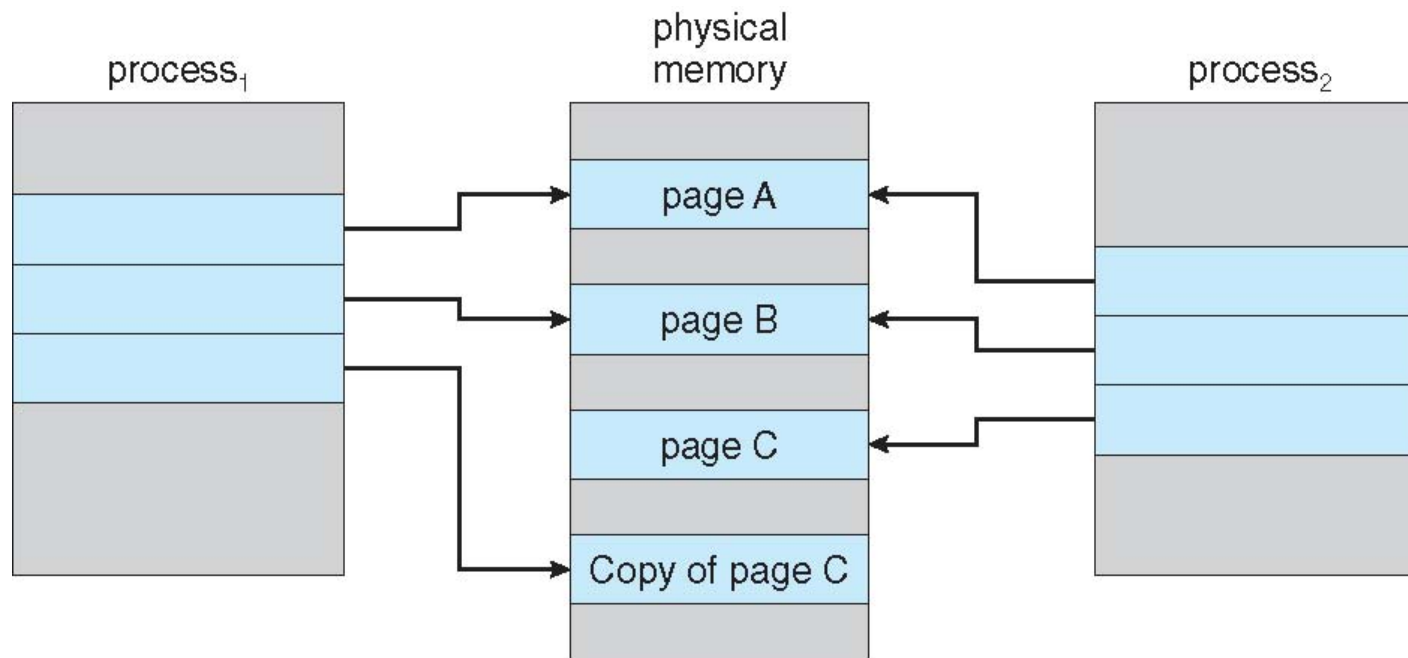    - Very efficient (page table is not copied)

# After fork()

# After Page C Modified



process₁

physical memory

process₂

page A

page B

page C

Copy of page C

# **What Happens if There is no Free Frame?**

- Used up by process pages

- Also in demand from the kernel, I/O buffers, etc

- How much to allocate to each?

- Page replacement – find some page in memory, but not really in use, page it out
    - Algorithm – terminate? swap out? replace the page?
    - Performance – want an algorithm which will result in minimum number of page faults

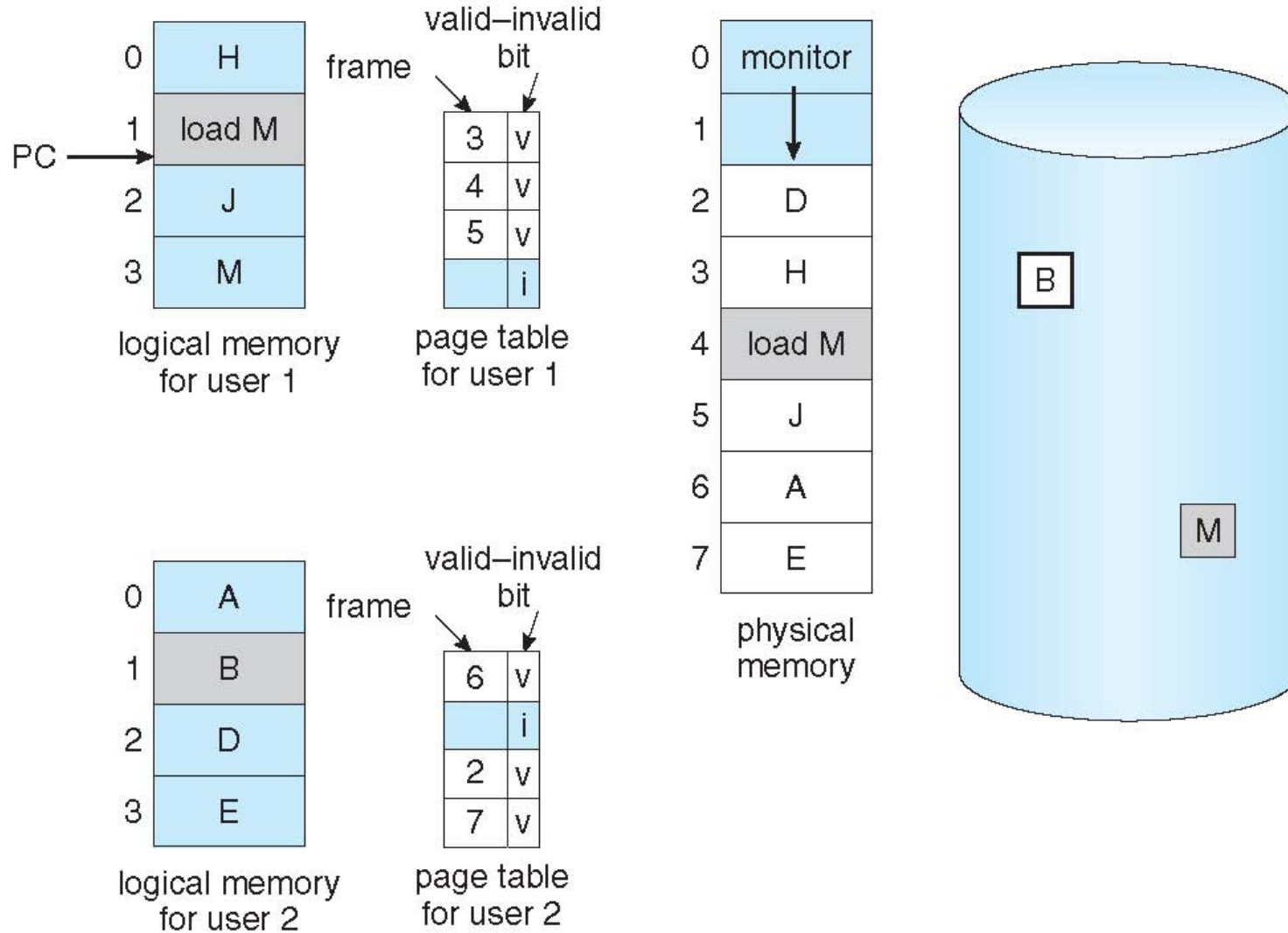- Same page may be brought into memory several times

# Page Replacement

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

- Use **modify (dirty) bit** to reduce overhead of page transfers – only modified pages are written to disk

- Page replacement completes separation between logical memory and physical memory – large virtual memory can be provided on a smaller physical memory

logical memory for user 1

page table for user 1

logical memory for user 2

page table for user 2

physical memory
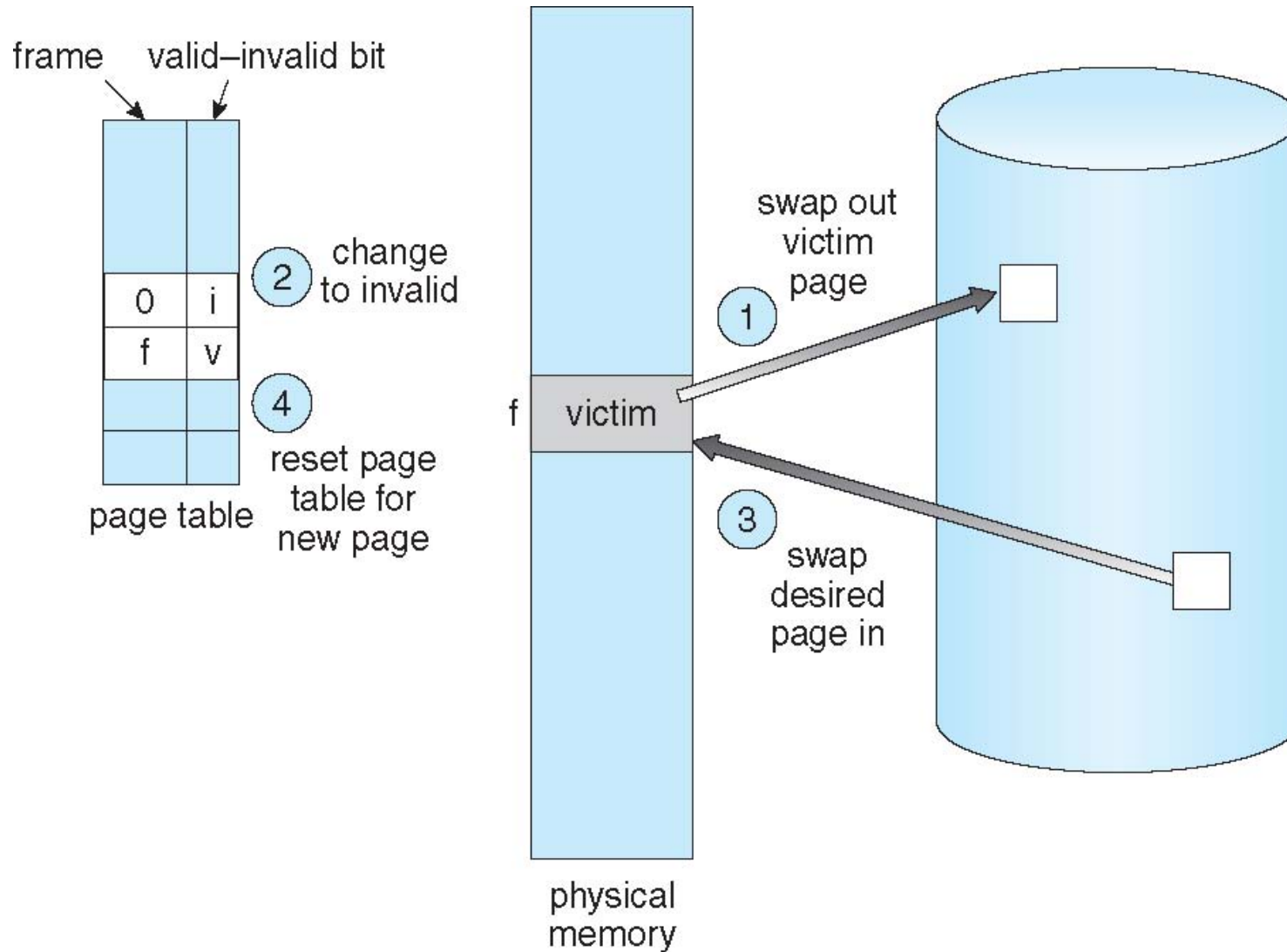
# Basic Page Replacement

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

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

page table

| 0 | i |
| f | v |

② change to invalid

④ reset page table for new page

f  victim

physical memory

① swap out victim page

③ swap desired page in

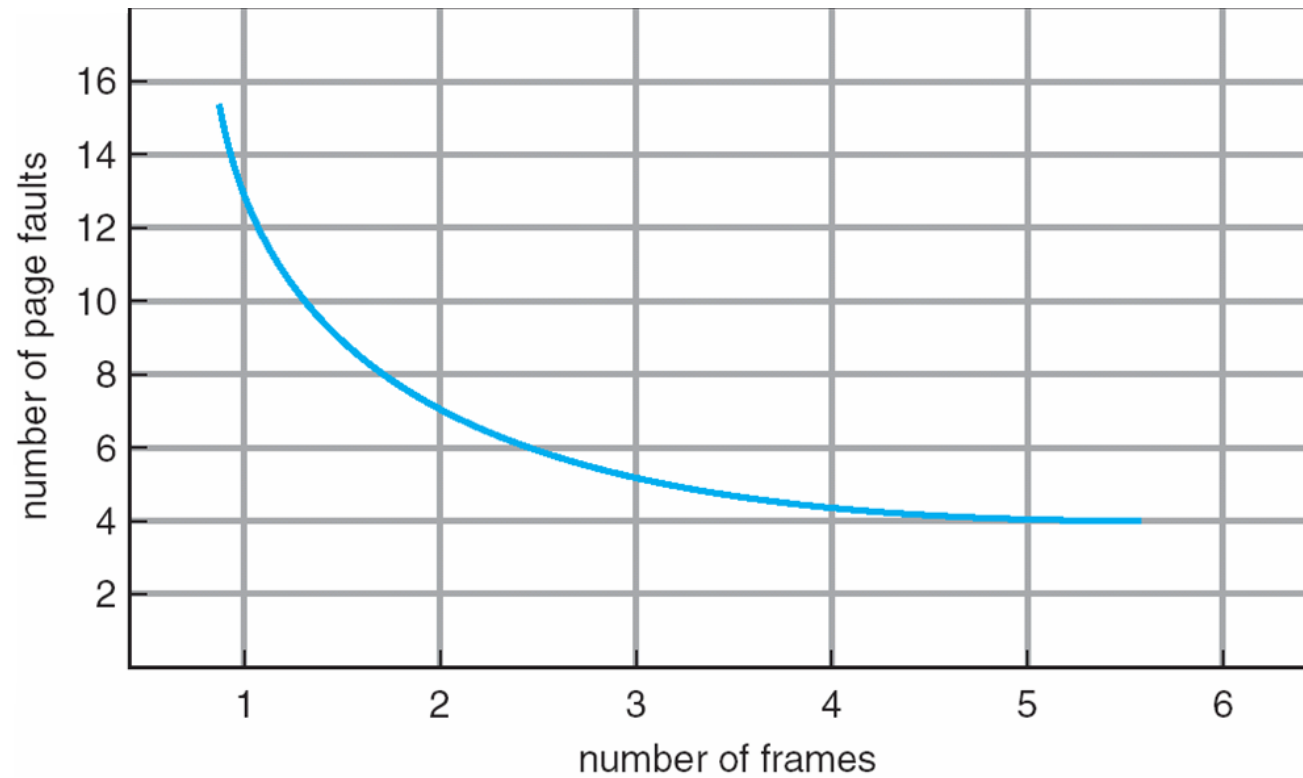# Page and Frame Replacement Algorithms

- **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

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

- In all our examples, memory has 3 frames and 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**

# 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)

# 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)

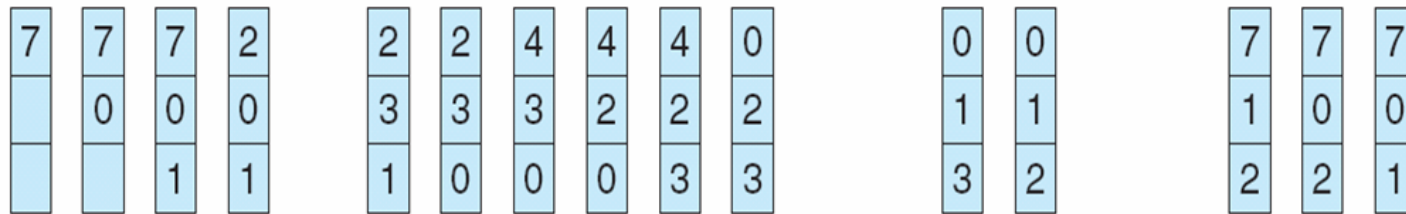reference string

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

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

page frames

- 15 page faults

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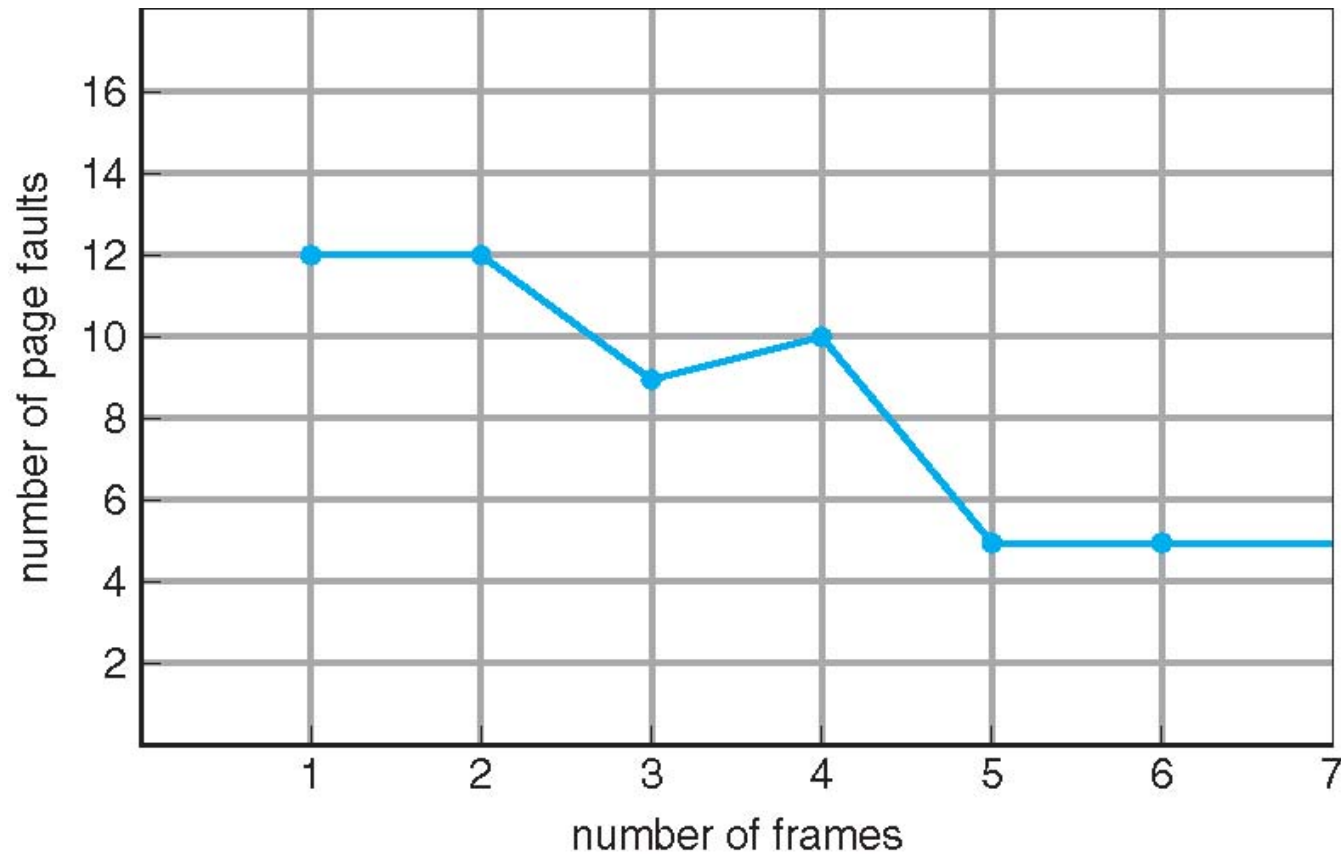  - Just use a FIFO queue

# Exercise

- Reference string: **1,2,3,4,1,2,5,1,2,3,4,5**
  - Evaluate FIFO replacement with **3 frames**
  - Evaluate FIFO replacement with **4 frames**

# Belady's Anomaly

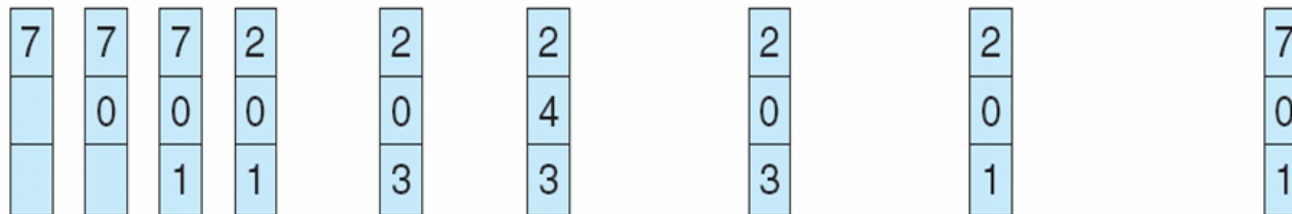■ Adding more frames can cause more page faults!

# Optimal Page Replacement Algorithm

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

- Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**

# Optimal Page Replacement Algorithm

■ Replace page that will not be used for longest period of time

■ Reference string: **7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**

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 Page Replacement Algorithm

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

- How do you know this?
  - Can't read the future

- Used for measuring how well your algorithm performs
  - **Lower bound**

# 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: **7,0,1,2,0,3,0,4,2,3,0,3,0,3,2,1,2,0,1,7,0,1**
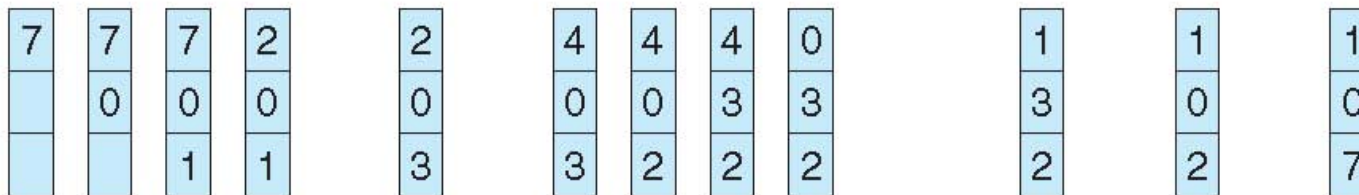
# Least Recently Used (LRU) Algorithm

- 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
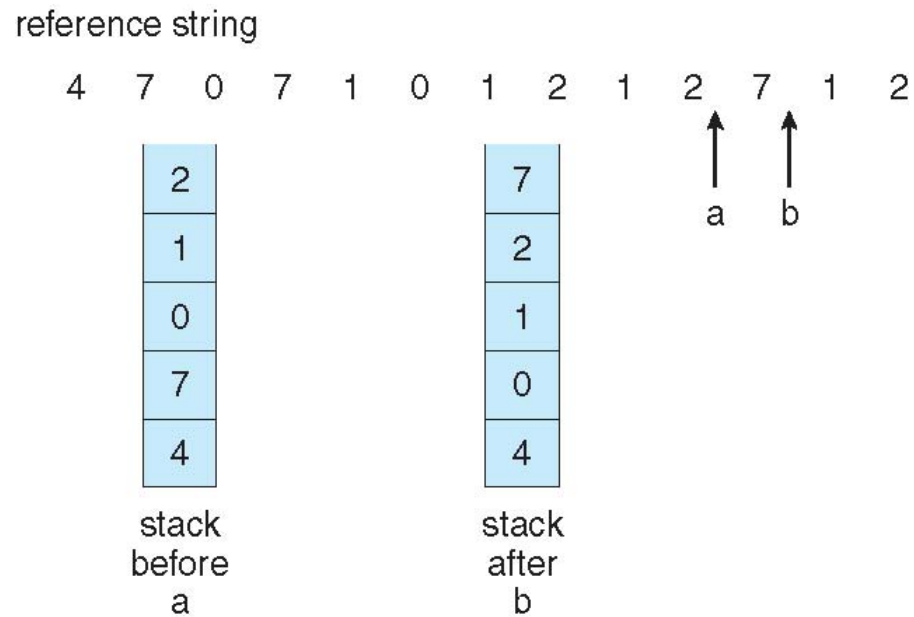- But how to implement?

# 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
  - Keep a stack of page numbers in a double link form:
  - Page referenced:
    - move it to the top
    - requires 6 pointers to be changed
  - No search for replacement
  - But each update more expensive
- LRU and OPT are cases of **stack algorithms** that don't have Belady's Anomaly

reference string

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

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

stack
before
a

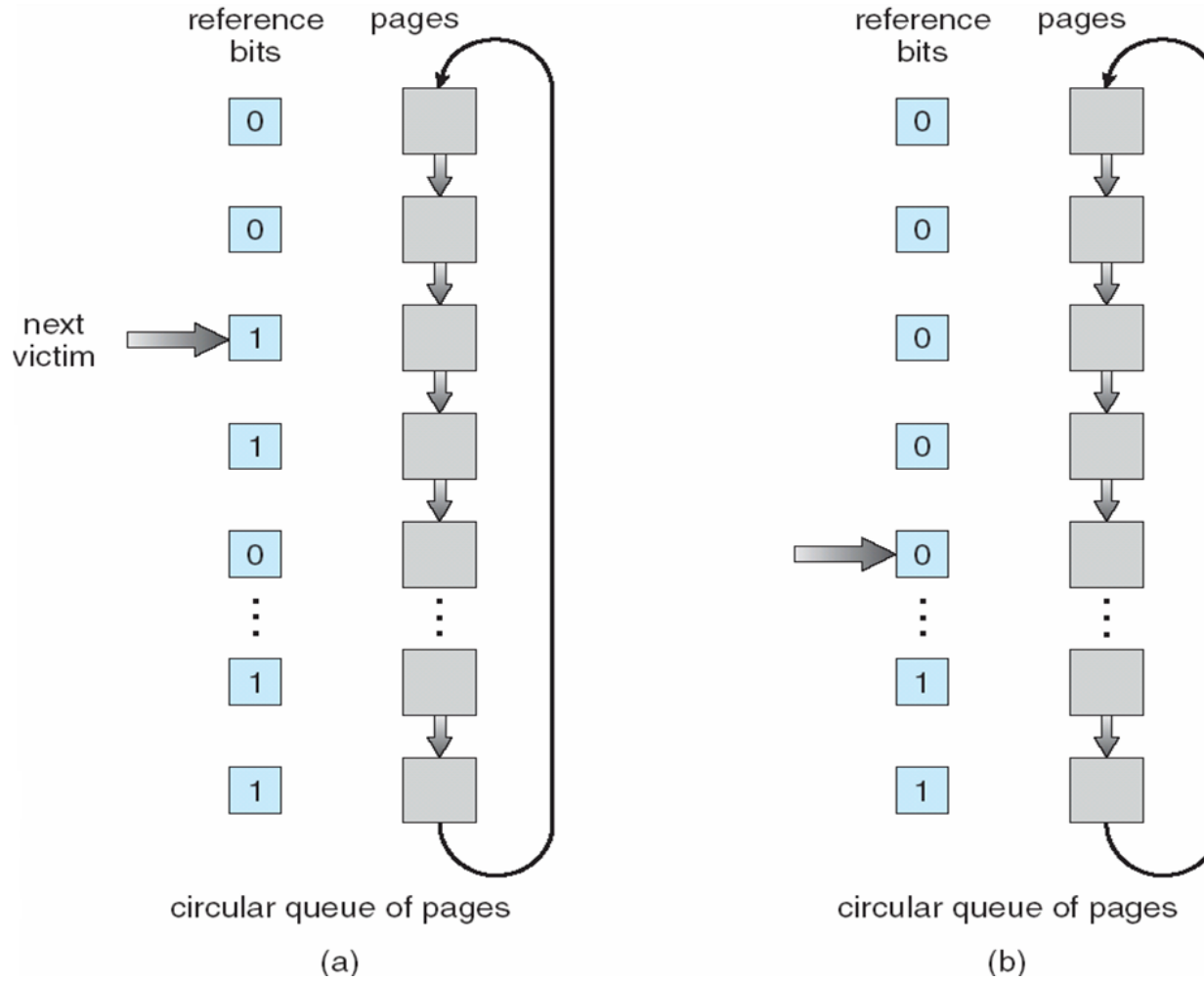| 7 |
|---|
| 2 |
| 1 |
| 0 |
| 4 |

stack
after
b

# 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**
    - 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) Page-Replacement Algorithm

# 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
- 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

# 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 $p_i = \dfrac{s_i}{S} \times m$

  $m = 64$

  $s_1 = 10$

  $s_2 = 127$

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

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

# 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
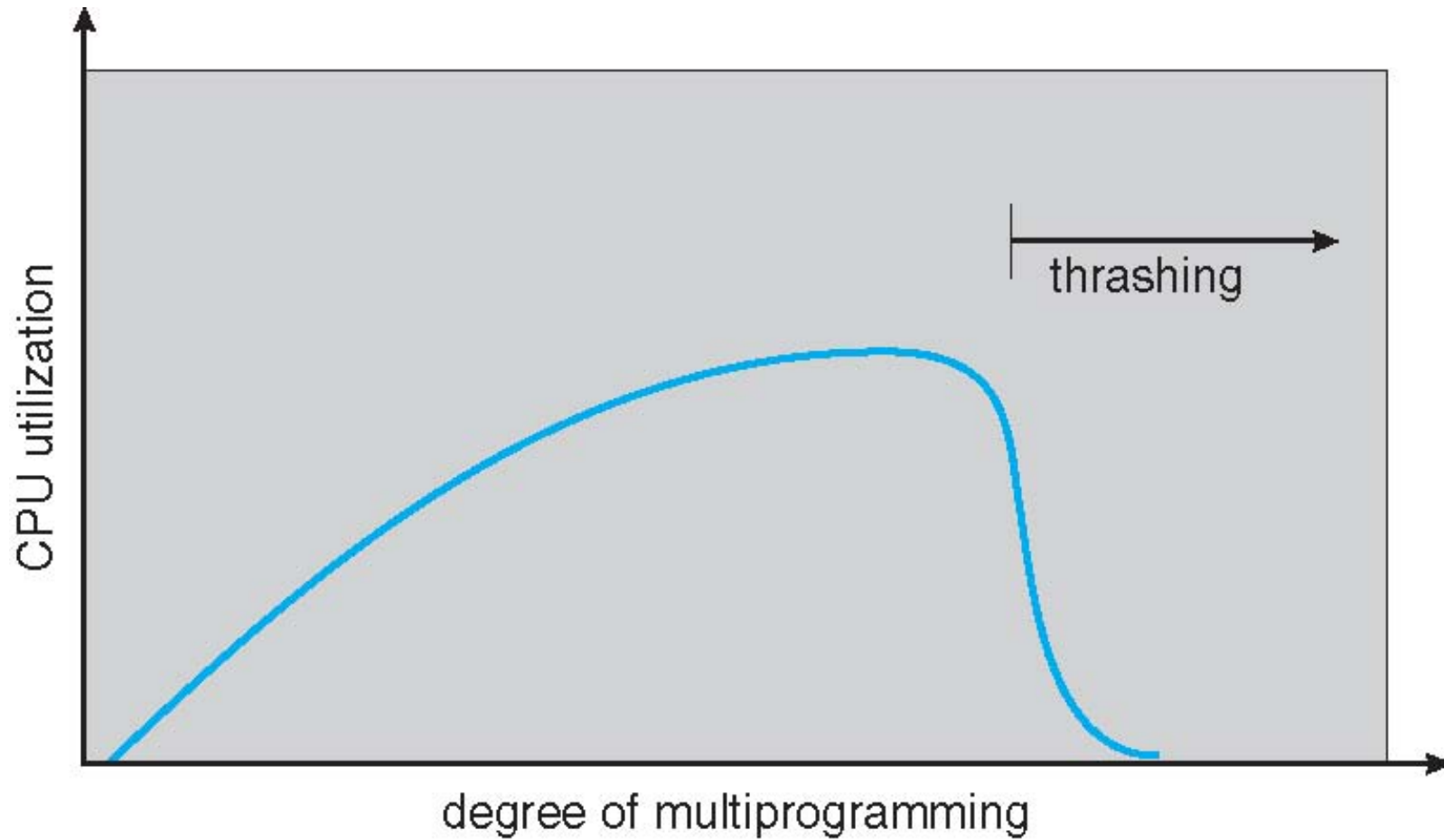  - But possibly underutilized memory

# Thrashing

■ 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.)



thrashing

CPU utilization

degree of multiprogramming

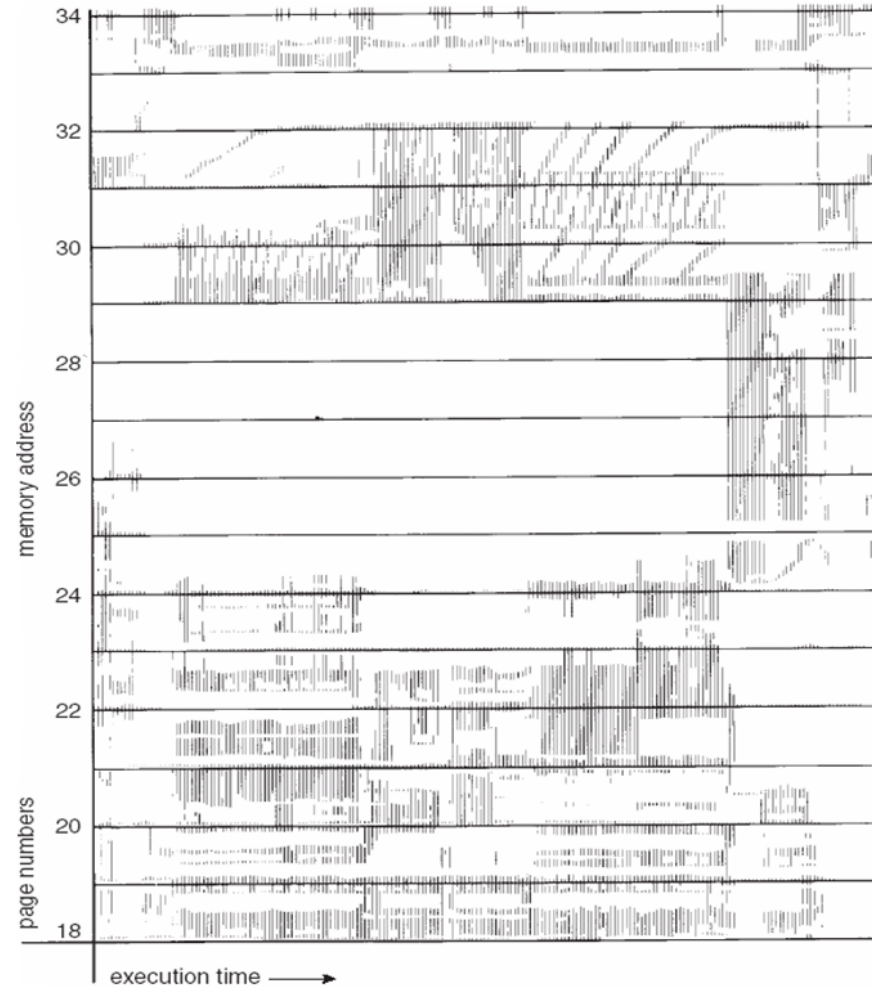# Demand Paging and Thrashing

- 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

- $\Delta \equiv$ 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

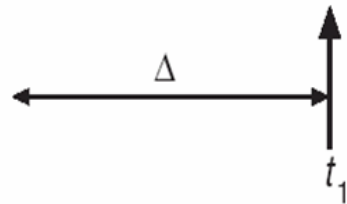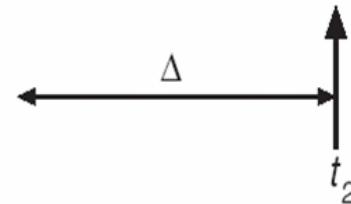- Policy if $D > m$, then suspend or swap out one of the processes

# Working-Set Model

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 . . .

$\Delta$

$\Delta$

$t_1$

$t_2$

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

$WS(t_2) = \{3,4\}$
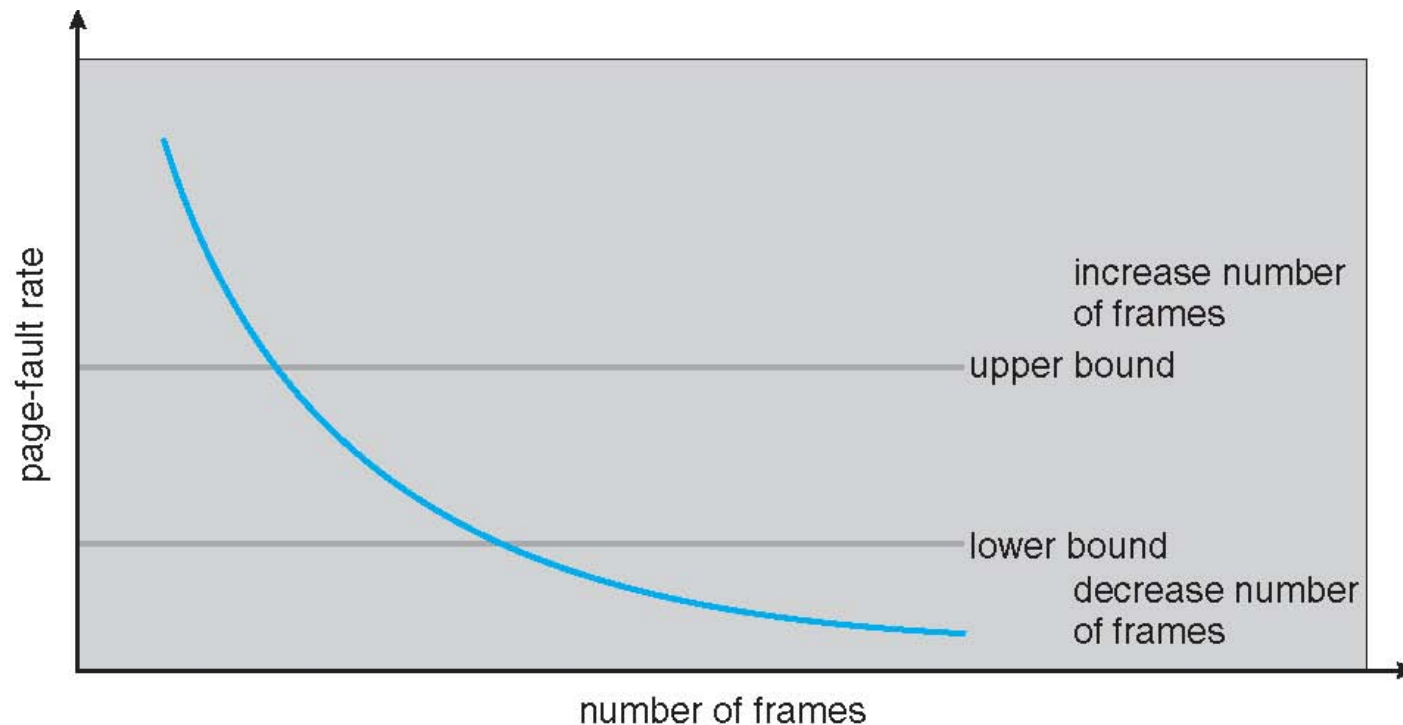
# Keeping Track of the Working Set

- 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?

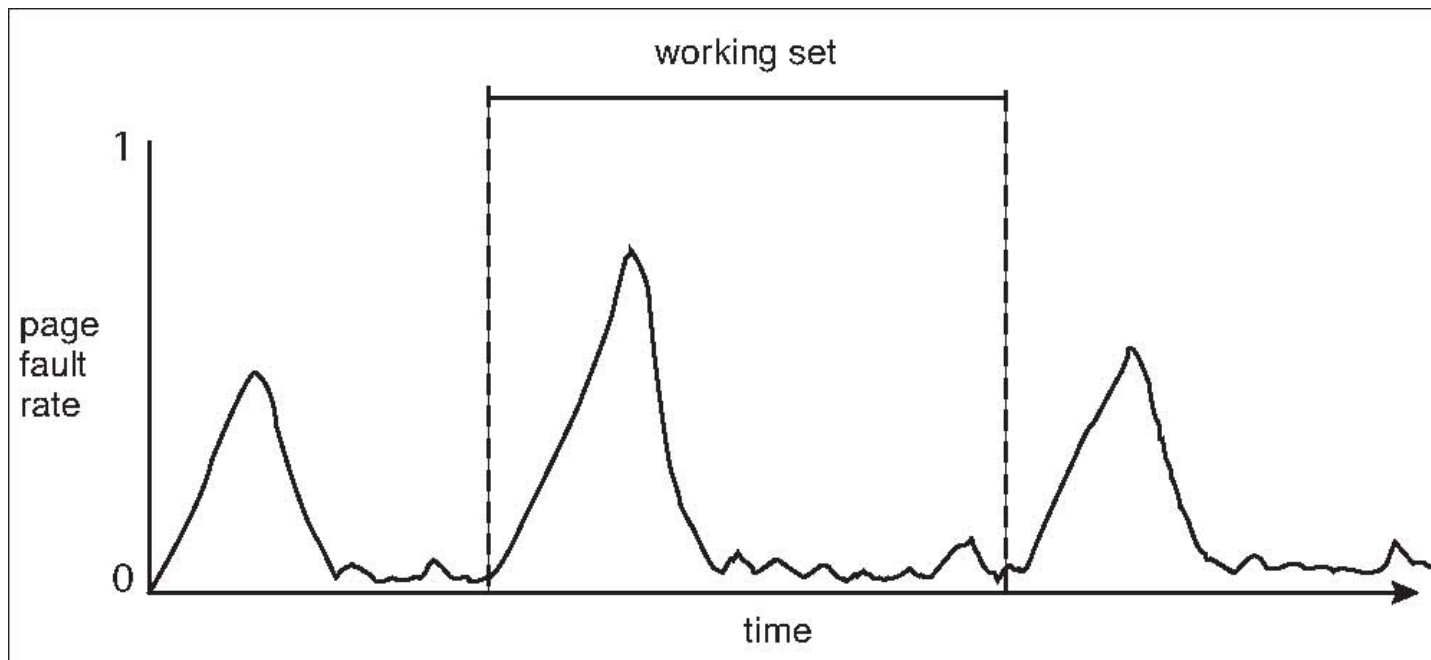- Improvement = 10 bits and interrupt every 1000 time units

# Page-Fault Frequency

- More direct approach than WSS

- Establish "acceptable" **page-fault frequency** 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

# Memory-Mapped Files

- Memory-mapped file I/O allows file I/O to be treated as routine memory access by **mapping** a disk block to a page in memory

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

- Simplifies and speeds file access by driving file I/O through memory rather than `read()` and `write()` system calls

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

- But when does written data make it to disk?
  - Periodically and / or at file `close()` time
  - For example, when the pager scans for dirty pages

# Memory-Mapped Files