# Chapter 8: Main Memory

# Memory Management – Main Memory

- Background

- Swapping

- Contiguous Allocation

- Paging

- Segmentation

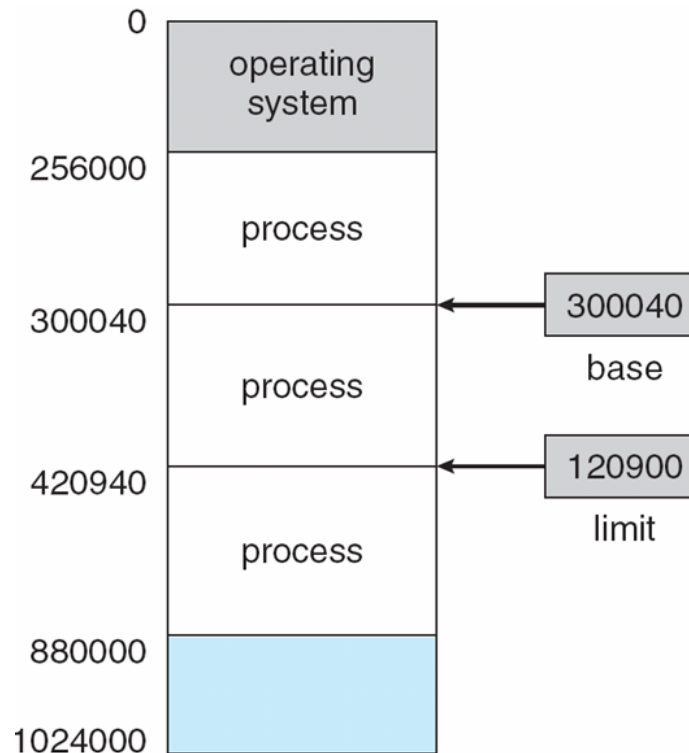- Segmentation with Paging

# Background

- Program must be brought (from disk) into memory and placed within a process for it to be run

- Main memory and registers are only storage CPU can access directly

- Memory unit only sees a stream of addresses + read requests, or address + data and write requests

- Register access in one CPU clock (or less)

- Main memory can take many cycles, causing a **stall**

- **Cache** sits between main memory and CPU registers

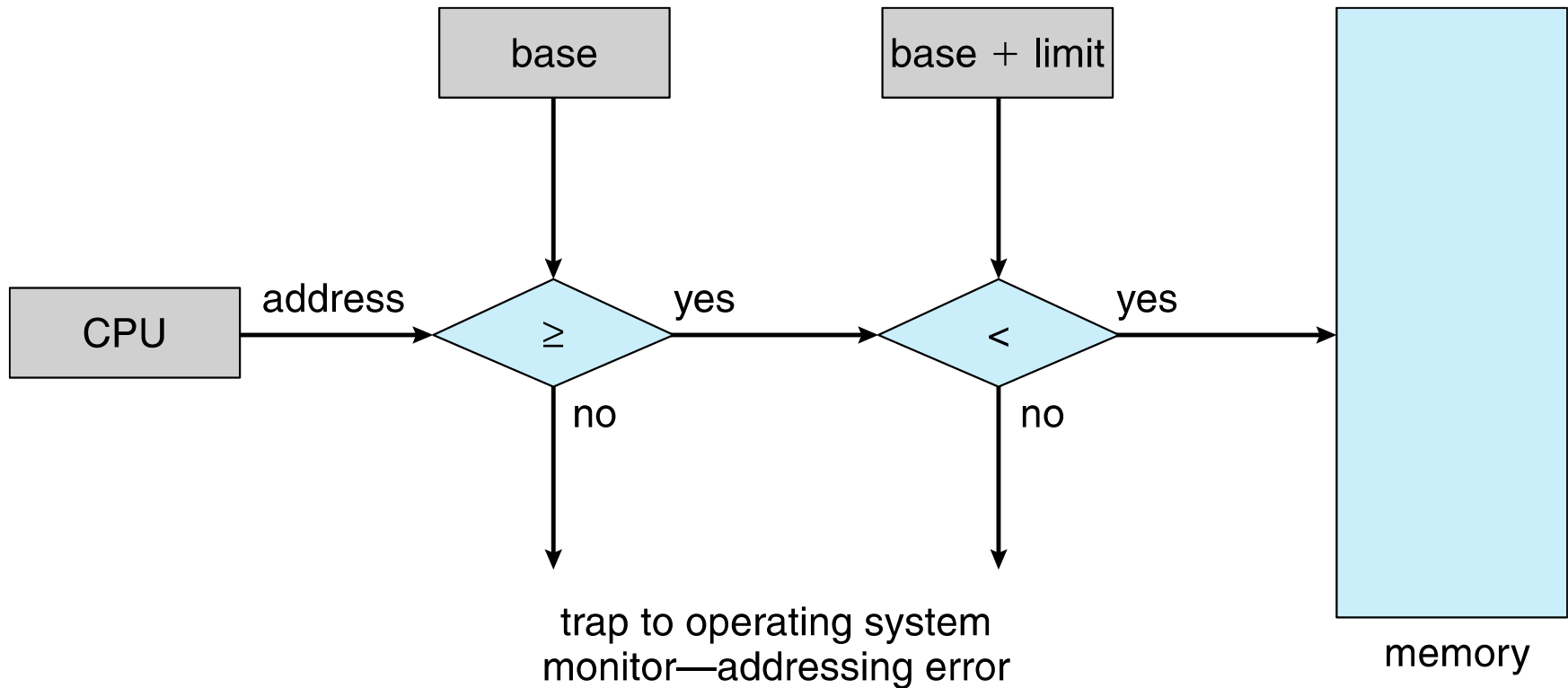- Protection of memory required to ensure correct operation

# Base and Limit Registers

- A pair of **base** and **limit registers** define the logical address space
- CPU must check every memory access generated in user mode to be sure it is between base and limit for that user

# Hardware Address Protection



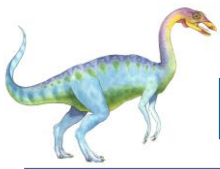trap to operating system
monitor—addressing error

memory

# Logical vs. Physical Address Space

- The concept of a logical address space that is bound to a separate **physical address space** is central to proper memory management

    - **Logical address** – generated by the CPU; also referred to as **virtual address**

    - **Physical address** – address seen by the memory unit

- Logical and physical addresses are the same in compile-time and load-time address-binding schemes; logical (virtual) and physical addresses differ in execution-time address-binding scheme

- **Logical address space** is the set of all logical addresses generated by a program

- **Physical address space** is the set of all physical addresses generated by a program
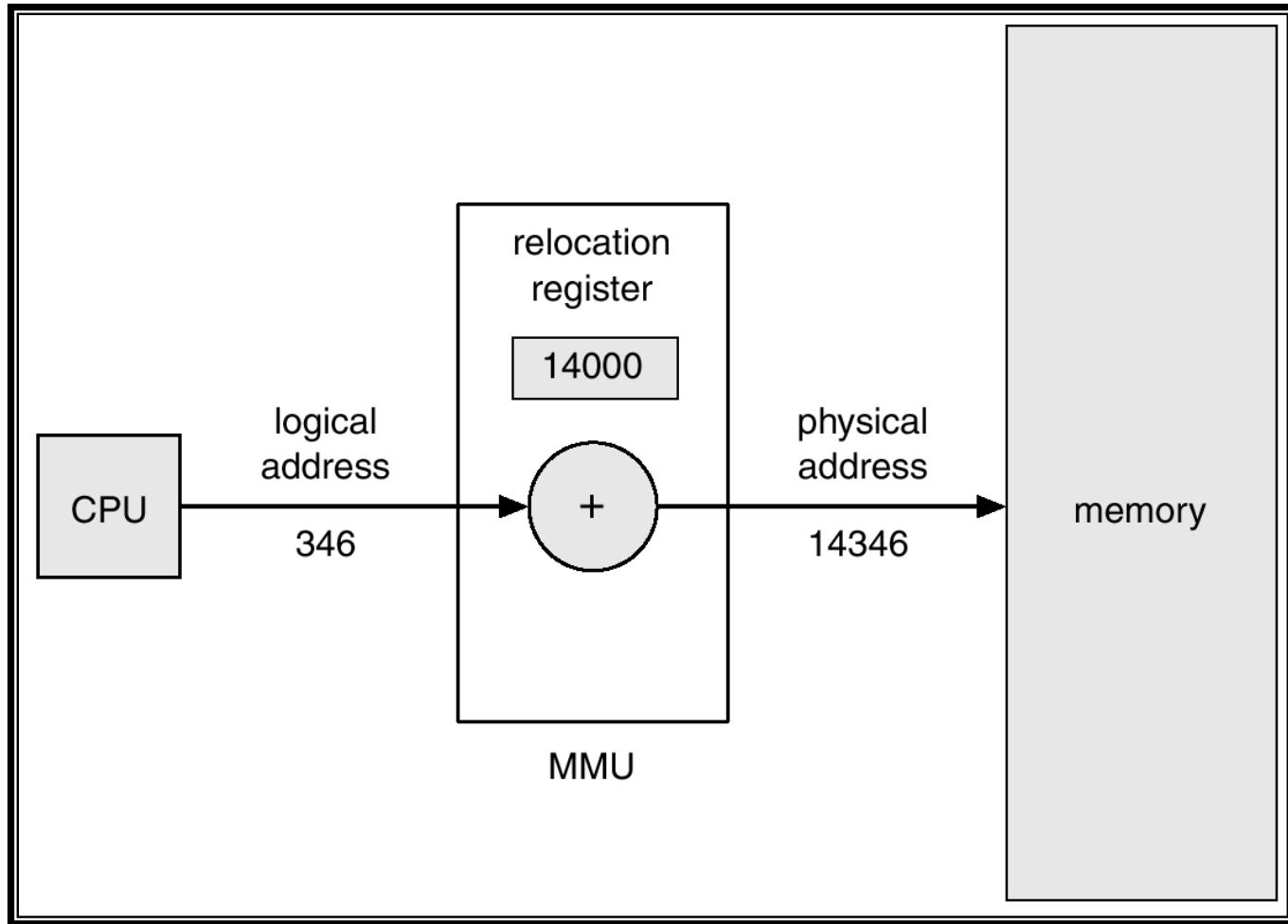
# Memory-Management Unit (MMU)

- Hardware device that maps virtual to physical address.

- In **MMU** scheme, the value in the relocation register is added to every address generated by a user process at the time it is sent to memory.

- The user program deals with *logical* addresses; it never sees the *real* physical addresses.
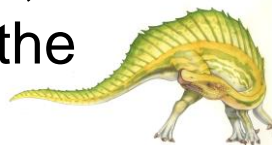
# Dynamic relocation using a relocation register

# Dynamic Linking and Shared Libraries

- With **Static linking** library modules get fully included in executable modules, wasting both disk space and main memory usage.

- With **dynamic linking**, however, only a stub is linked into the executable module, containing references to the actual library module linked in at run time.

  - Saves disk space.

  - it does not modify the code while it runs, making it safe to re-enter.

  - An added benefit of **dynamically linked libraries** (**DLLs**), involves easy upgrades and updates.

  - In practice, the first time a program calls a DLL routine, the stub will recognize the fact and will replace itself with the actual routine from the DLL library.
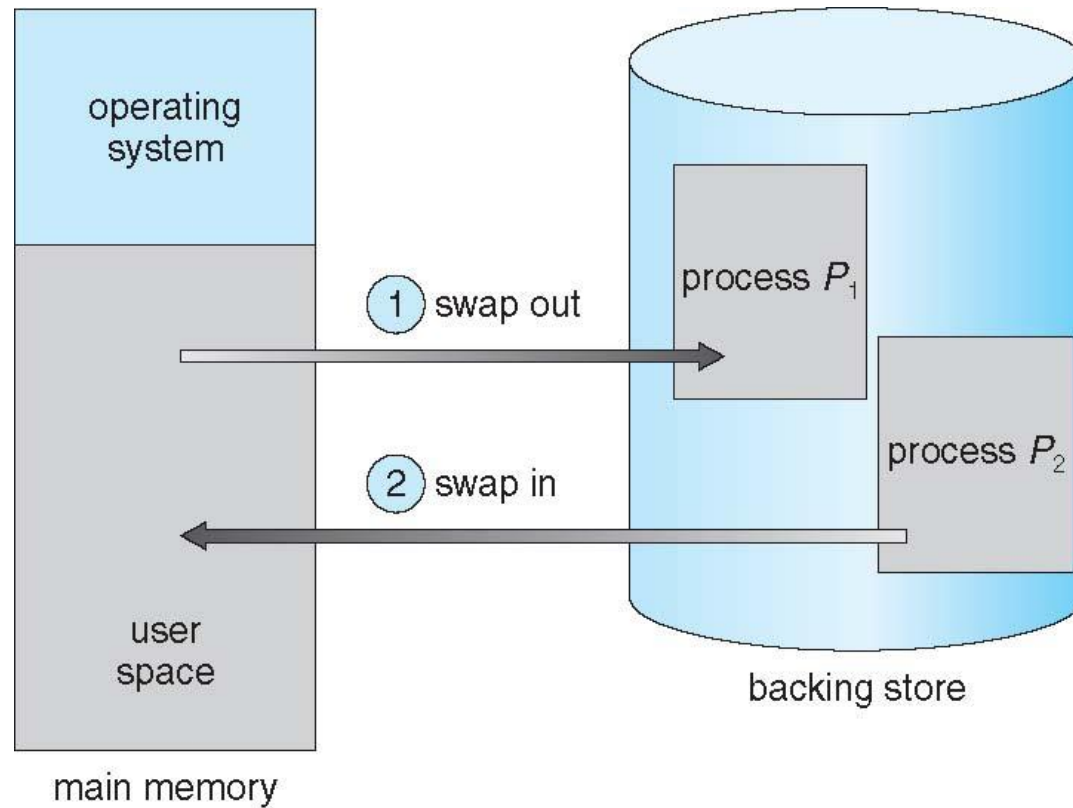
# Swapping

- A process can be *swapped* temporarily out of memory to a *backing store*, and then brought back into memory for continued execution.

- Backing store – fast disk large enough to accommodate copies of all memory images for all users; must provide direct access to these memory images.

- *Roll out, roll in* – swapping variant used for priority-based scheduling algorithms; lower-priority process is swapped out so higher-priority process can be loaded and executed.

- Major part of swap time is transfer time; total transfer time is directly proportional to the *amount* of memory swapped.

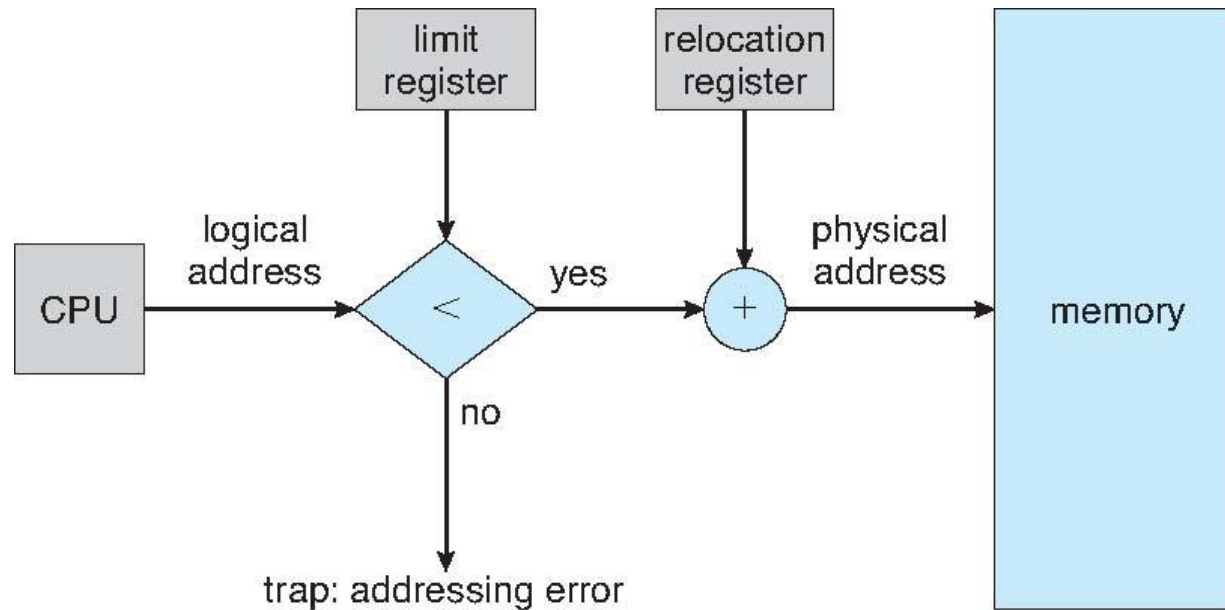- Modified versions of swapping are found on many systems, i.e., UNIX, Linux, and Windows.
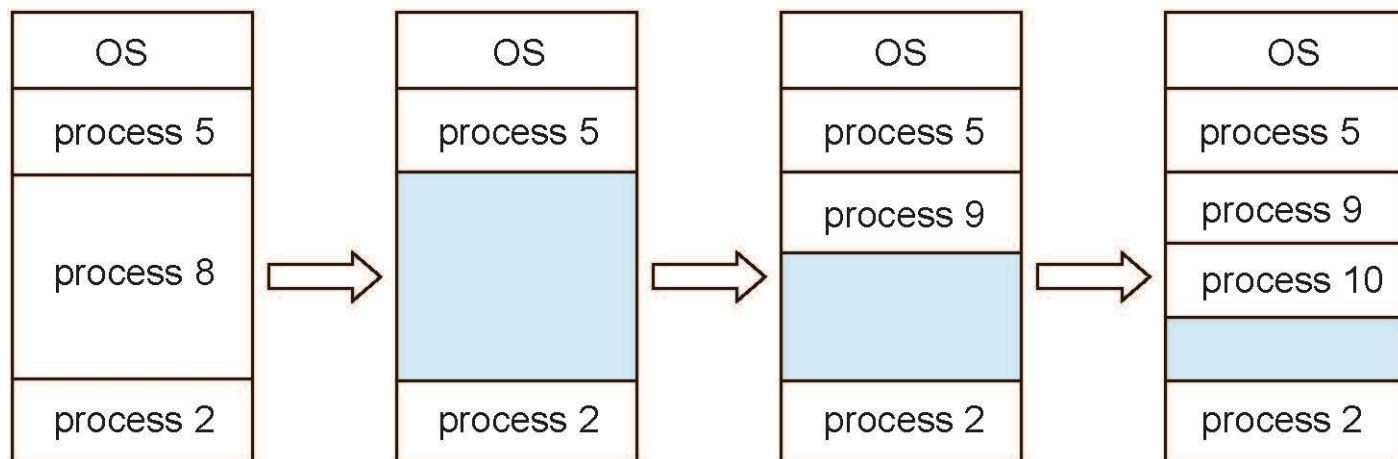
# Schematic View of Swapping

# Multiple-partition allocation

- **Multiple-partition allocation**
    - Degree of multiprogramming limited by number of partitions
    - **Variable-partition** sizes for efficiency (sized to a given process' needs)
    - **Hole** – block of available memory; holes of various size are scattered throughout memory
    - When a process arrives, it is allocated memory from a hole large enough to accommodate it
    - Process exiting frees its partition, adjacent free partitions combined
    - Operating system maintains information about:
      a) allocated partitions    b) free partitions (hole)

| OS |
|----|
| process 5 |
| process 8 |
| process 2 |

→

| OS |
|----|
| process 5 |
|  |
| process 2 |

→

| OS |
|----|
| process 5 |
| process 9 |
|  |
| process 2 |

→

| OS |
|----|
| process 5 |
| process 9 |
| process 10 |
|  |
| process 2 |

# Dynamic Storage-Allocation Problem

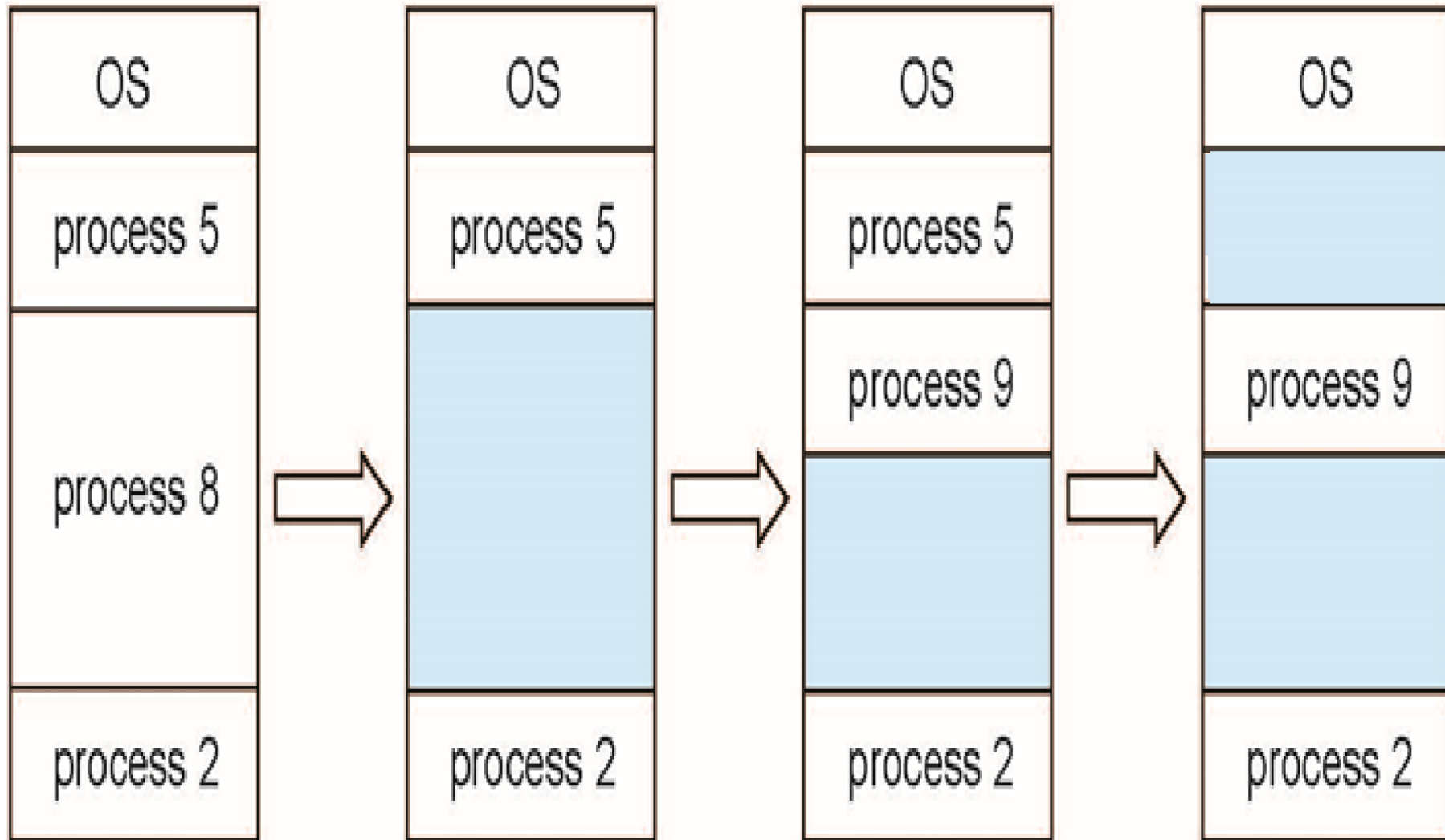How to satisfy a request of size $n$ from a list of free holes?

- **First-fit**:  Allocate the *first* hole that is big enough

- **Best-fit**:  Allocate the *smallest* hole that is big enough; must search entire list, unless ordered by size
  - Produces the smallest leftover hole

- **Worst-fit**:  Allocate the *largest* hole; must also search entire list
  - Produces the largest leftover hole

First-fit and best-fit better than worst-fit in terms of speed and storage utilization

# Fragmentation

# Paging

- Logical address space of a process can be noncontiguous; process is allocated physical memory whenever the latter is available.

- Divide physical memory into fixed-sized blocks called **frames** (size is power of 2, between 512 bytes and 8192 bytes).

- Divide logical memory into blocks of same size called **pages**.

- Keep track of all free frames.

- To run a program of size *n* pages, need to find *n* free frames and load program.

- Set up a page table to translate logical to physical addresses.

- Internal fragmentation.

# Address Translation Architecture

# Paging Example

# Paging Example



logical memory

page table

physical memory

# Free Frames



Before allocation                    After allocation

# 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* (PRLR) 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)*

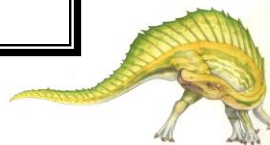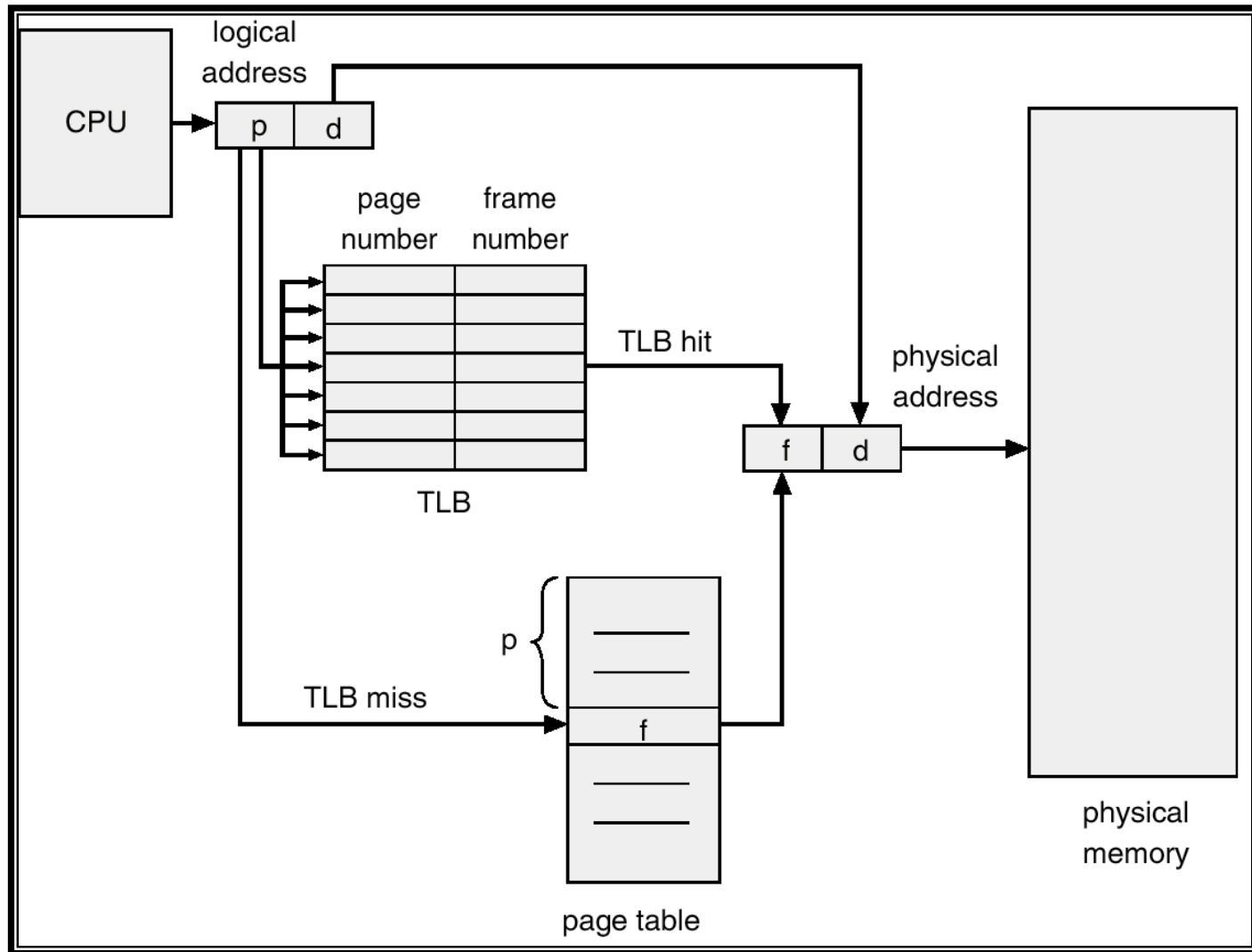# Associative Memory

- Associative memory – parallel search

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

Address translation (A´, A´´)

- If A´ is in associative register, get frame # out.
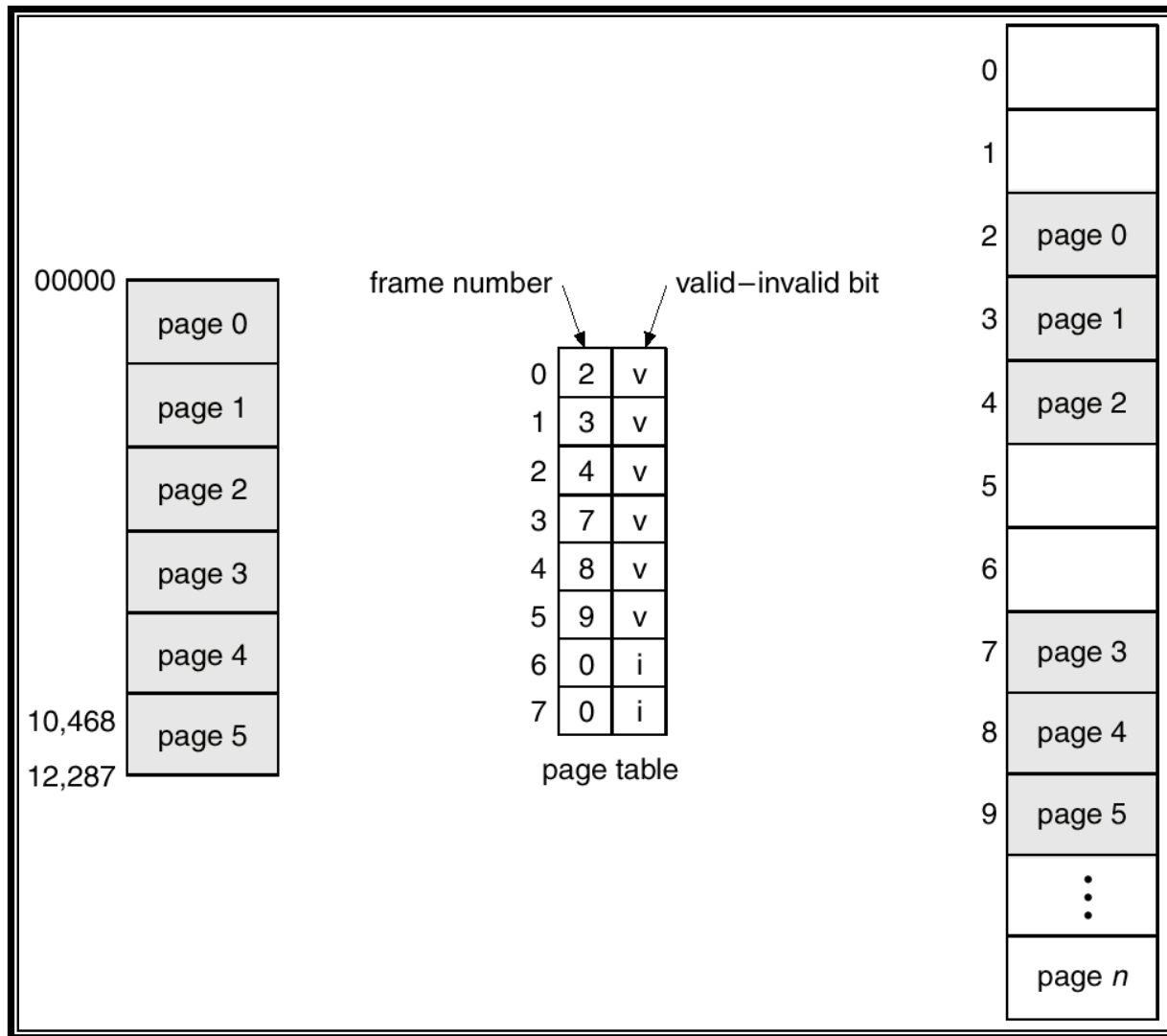- Otherwise get frame # from page table in memory

# Paging Hardware With TLB

# Memory Protection

- Memory protection implemented by associating protection bit with each frame.

- *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.

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

# Page Table Structure

- 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.

# 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 |
|---|---|---|
| $p_i$ | $p_2$ | $d$ |
| 10 | 10 | 12 |

where $p_i$ is an index into the outer page table, and $p_2$ is the displacement within the page of the outer page table.

# Two-Level Page-Table Scheme

# Address-Translation Scheme

- Address-translation scheme for a two-level 32-bit paging architecture

# Shared Pages

- **Shared code**

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

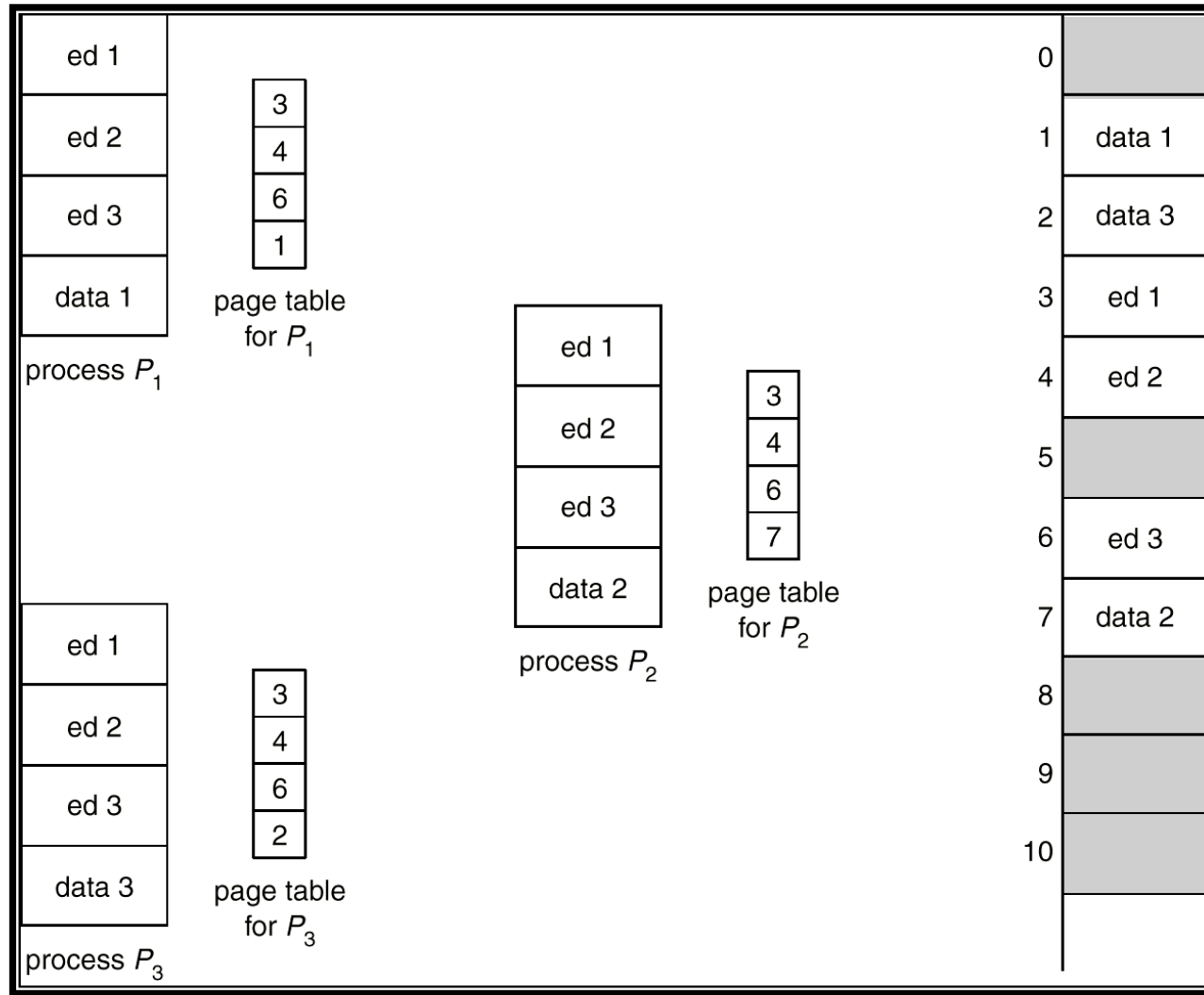  - Shared code must appear in same location in the logical address space of all processes.

- **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.

# Segmentation

■ Memory-management scheme that supports user view of memory.

■ A program is a collection of segments.  A segment is a logical unit such as:

main program,

procedure,

function,

method,

object,

local variables, global variables,

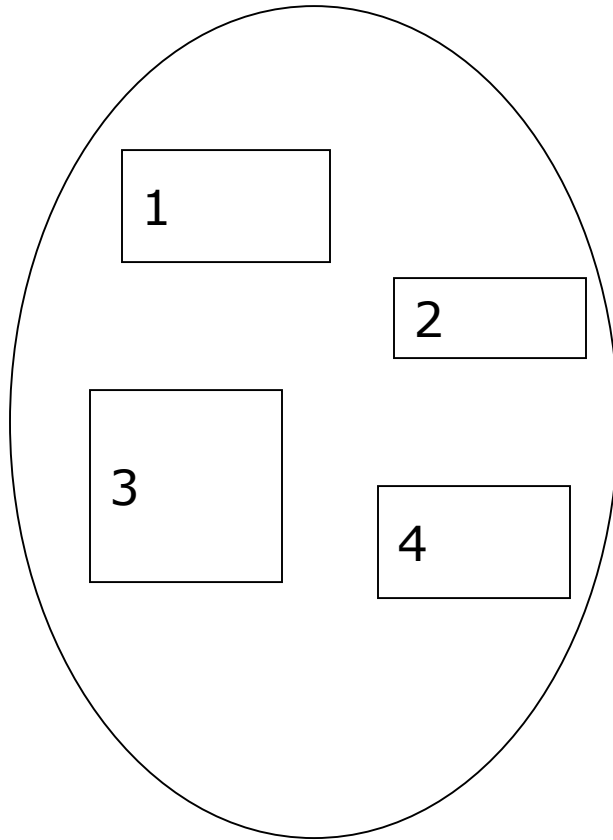common block,

stack,

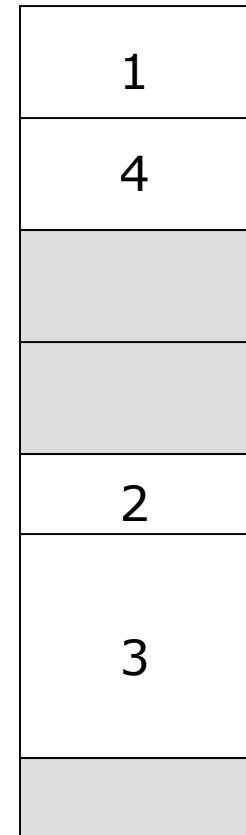symbol table, arrays

# User's View of a Program

# Logical View of Segmentation



user space

physical memory space

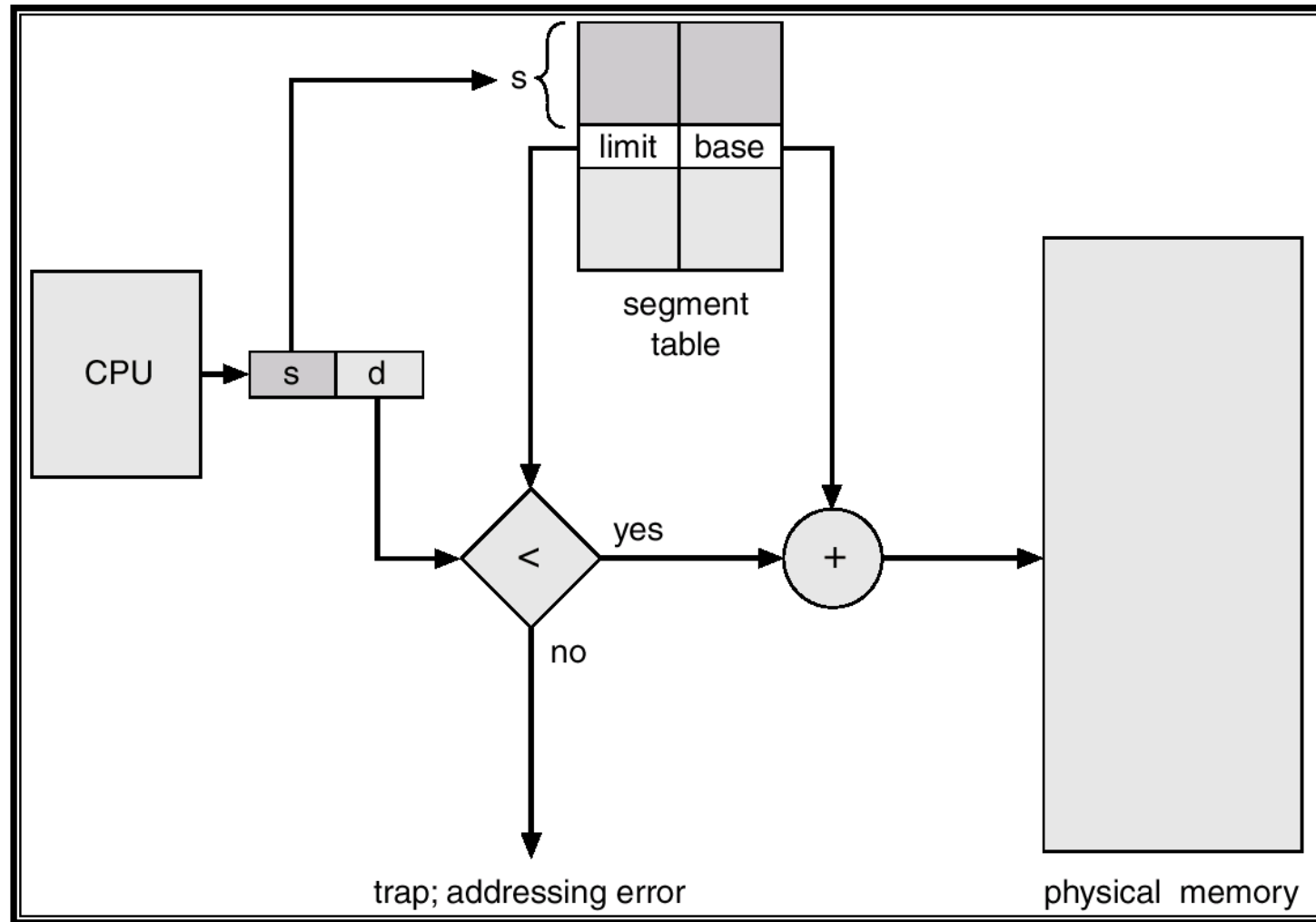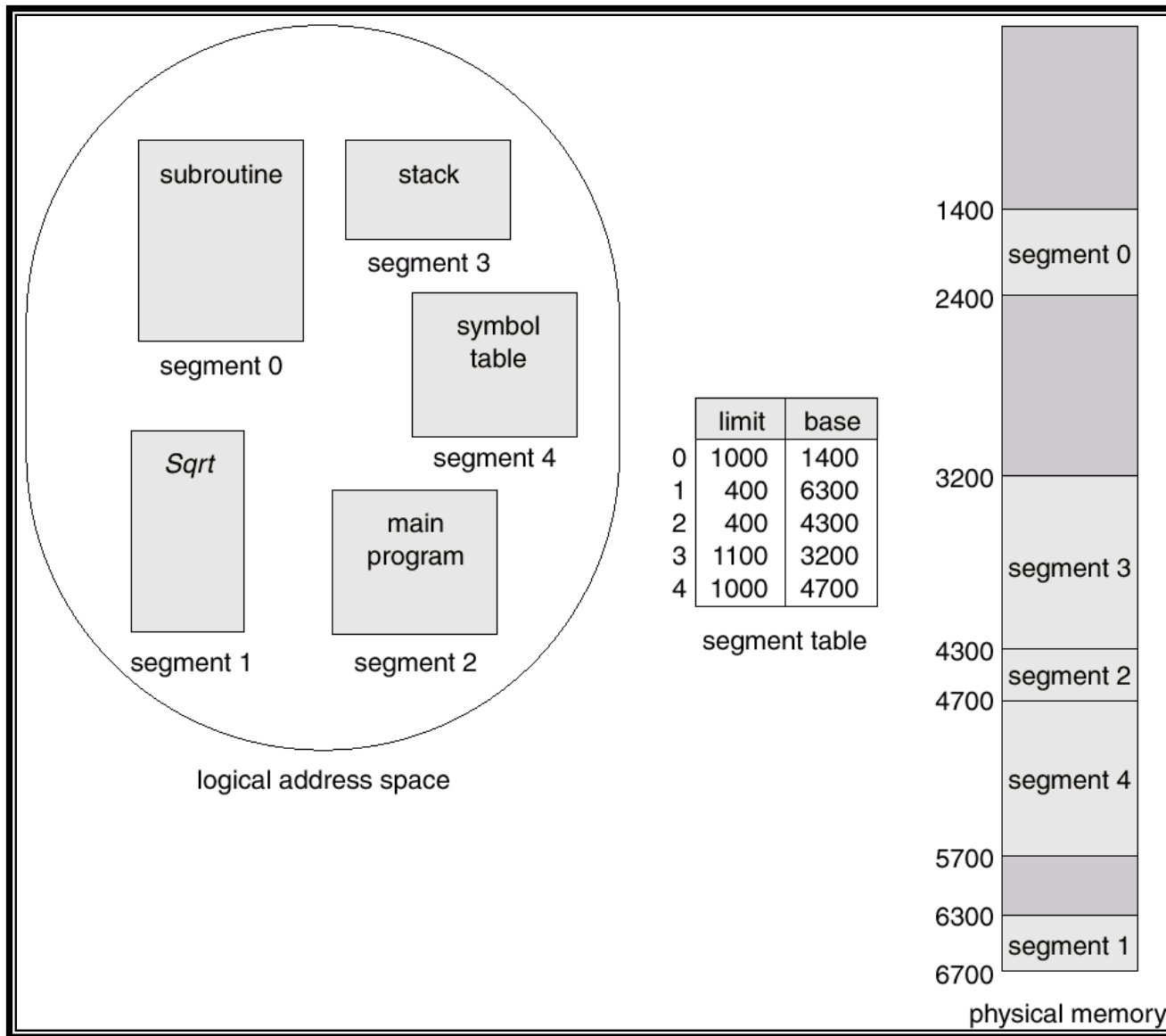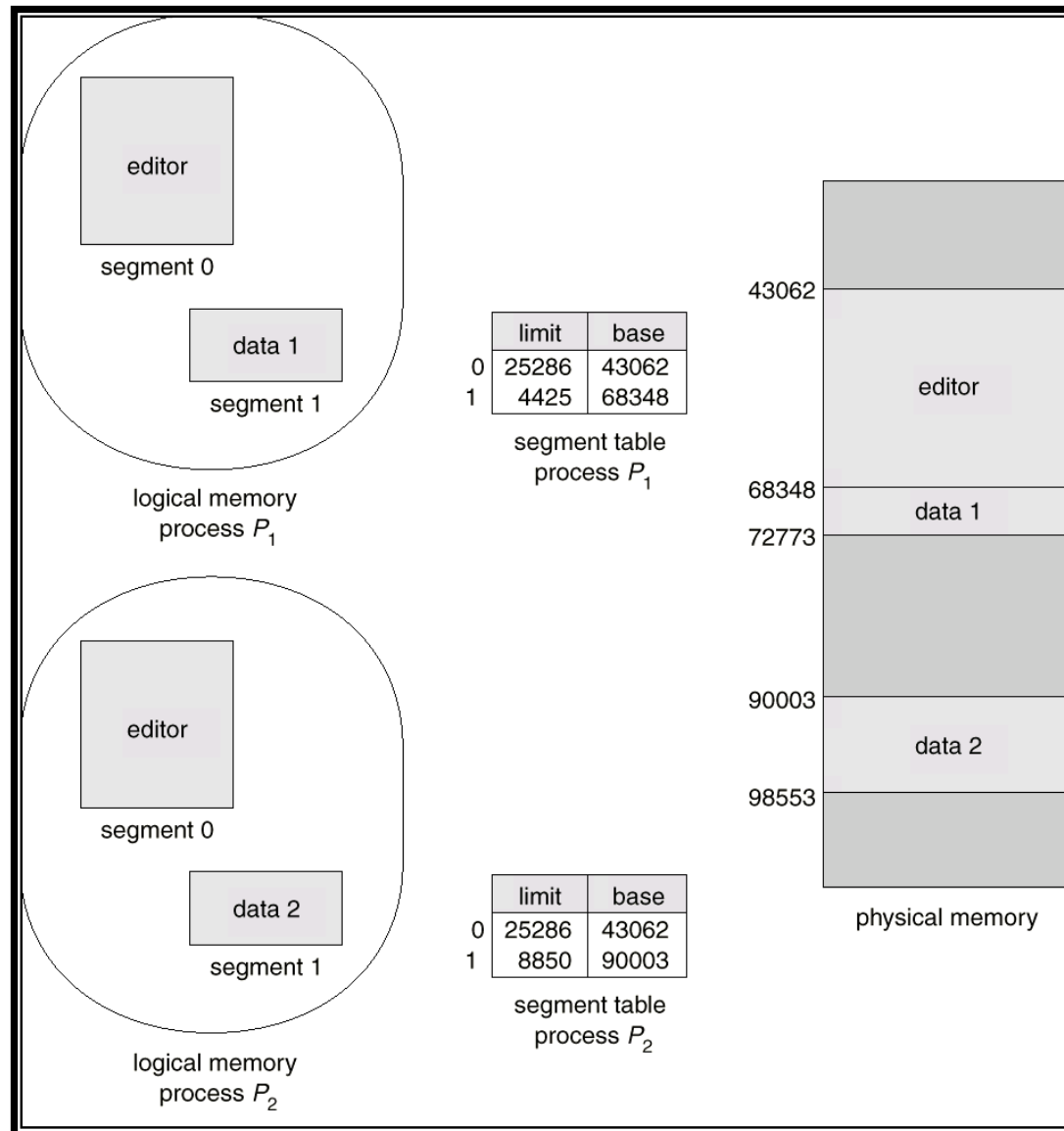# Segmentation Hardware

# Example of Segmentation

# Sharing of Segments

# End of Chapter 8