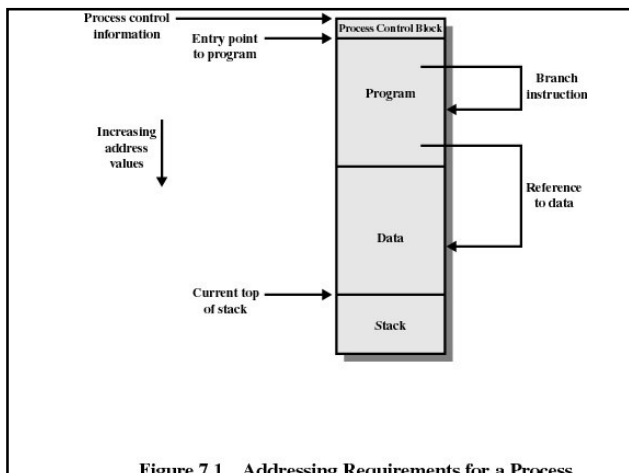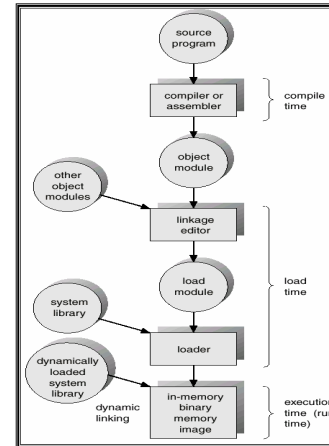# Memory Management

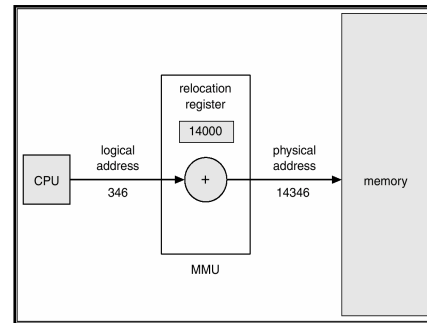Figure 7.1   Addressing Requirements for a Process

# Memory Management

- The challenge is:
  - Sub-divide the memory to accommodate multiple processes.

- Memory needs to be allocated efficiently to pack as many processes into memory as possible.

## Logical vs Physical Addresses

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

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

## Dynamic Relocation



## 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**.

## Memory Relocation

- **Relocation**
  - Programmer does not know where the program will be placed in memory when it is executed.
  - While the program is executing, it may be swapped to disk and returned to main memory at a different location (**relocated**).
  - Memory references must be translated in the code to actual physical memory address.
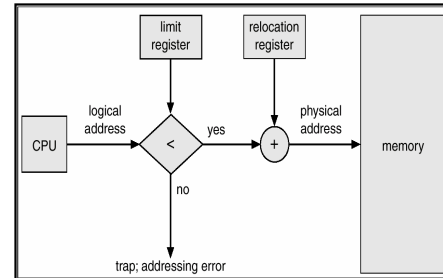
# Memory Protection

- **Protection**
  - Processes should not be able to reference memory locations in another process without permission.
  - Impossible to check absolute addresses in programs since the program could be relocated.
  - Must be checked during execution
    - Operating system cannot anticipate all of the memory references a program will make.

**Mechanisms that support relocation also support protection (hardware)**

# Memory Protection



Relocation-register: used to protect user processes from each other. Relocation register contains the value of the smallest physical address;
Limit register contains range of logical addresses – each logical address must be less than the limit register.

# Memory Sharing

- **Sharing**
  - Allow several processes to access the same portion of memory.
  - Better to allow each process access to the same copy of the program rather than have their own separate copy.

# Logical Organization

- **Logical Organization**
  - Programs are written in modules.
  - Modules can be written and compiled independently.
  - Different degrees of protection given to modules (read-only, execute-only).
  - Share modules.

## Physical Organization

- **Physical Organization**
  - Memory available for a process may be insufficient
    - Overlaying allows various modules to be assigned the same region of memory
  - Programmer does not know how much space will be available.

## Dynamic Loading

- Routine is not loaded until it is called.
- Better memory-space utilization; unused routine is never loaded.
- Useful when large amounts of code are needed to handle infrequently occurring cases.
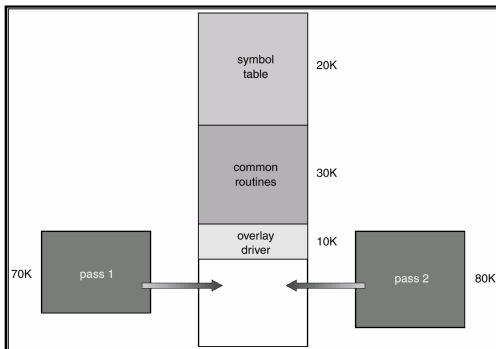
## Dynamic Linking

- Linking postponed until execution time.
- Small piece of code, *stub*, used to locate the appropriate memory-resident library routine.
- Stub replaces itself with the address of the routine, and executes the routine.
- Operating system needs to check if routine is in processes' memory address.
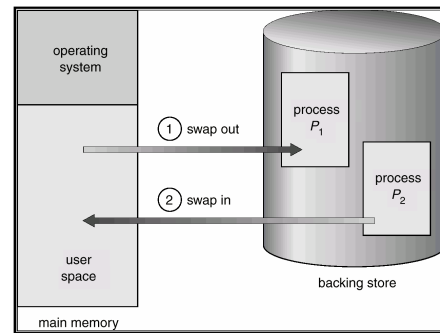- Dynamic linking is particularly useful for libraries.

## Memory Overlay

- Keep in memory only those instructions and data that are needed at any given time.
- Needed when process is larger than amount of memory allocated to it.
- Implemented by user, no special support needed from operating system.
- Programming design of overlay structure is complex

4

## Memory Overlay

symbol
table — 20K

common
routines — 30K

overlay
driver — 10K

70K — pass 1

pass 2 — 80K

## Swapping

operating
system

1 swap out

process $P_1$

2 swap in

process $P_2$

user
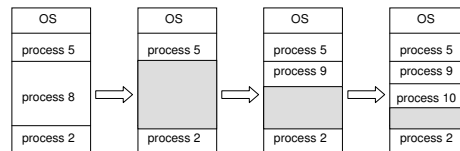space

backing store

main memory

## Swapping

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

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

- Example:
  disk: transfer rate       5Mb/sec
  memory to swap:       1 Mb
  Time it takes to swap out:     1/5 sec = 200 msec
  Average disk latency:       8 msec
  Time to swap-out/swap-in:    416 msec

## Memory Allocation

- Multiple-partition allocation
  - *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.
  - Operating system maintains information about:
    a) allocated partitions    b) free partitions (hole)

| OS | | OS | | OS | | OS |
|---|---|---|---|---|---|---|
| process 5 | | process 5 | | process 5 | | process 5 |
| | | | | process 9 | | process 9 |
| process 8 | ⇒ | | ⇒ | | ⇒ | process 10 |
| | | | | | | |
| process 2 | | process 2 | | process 2 | | process 2 |

# Memory Management Techniques

1- Fixed Partitioning

2- Dynamic Partitioning

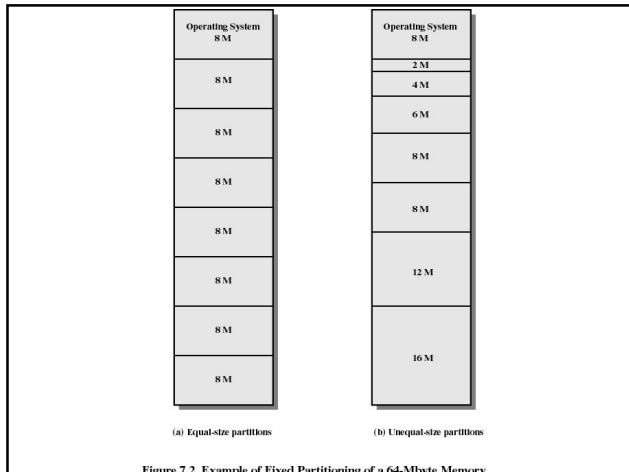3- Buddy System

4- Paging

5- Segmentation

# 1- Fixed Partitioning

# 1-Fixed Partitioning

- **Equal-size partitions**
  - any process whose size is less than or equal to the partition size can be loaded into an available partition.
  - if all partitions are full, the operating system can swap a process out of a partition.
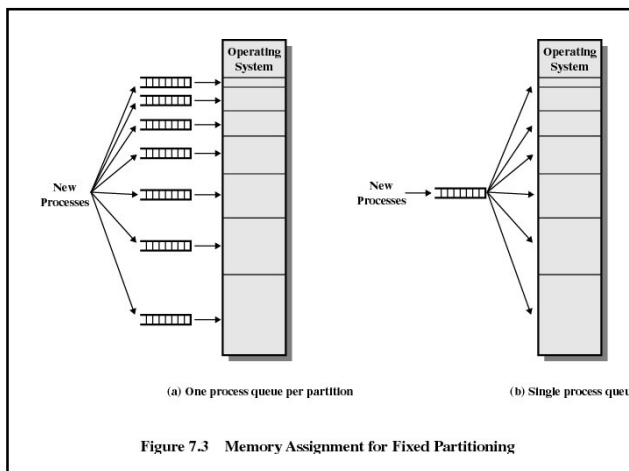  - a program may not fit in a partition. The programmer must design the program with overlays.

# Fixed Partitioning

- Main memory use is inefficient. Any program, no matter how small, occupies an entire partition. This is called **internal fragmentation.**

| Operating System 8 M | Operating System 8 M |
| 8 M | 2 M |
| | 4 M |
| 8 M | 6 M |
| 8 M | 8 M |
| 8 M | 8 M |
| 8 M | 12 M |
| 8 M | 16 M |
| 8 M | |
| (a) Equal-size partitions | (b) Unequal-size partitions |

Figure 7.2  Example of Fixed Partitioning of a 64-Mbyte Memory

## Placement Algorithm with Partitions

- **Equal-size partitions**
  - because all partitions are of equal size, it does not matter which partition is used
- **Unequal-size partitions**
  - can assign each process to the smallest partition within which it will fit
  - queue for each partition
  - processes are assigned in such a way as to minimize wasted memory within a partition



(a) One process queue per partition    (b) Single process queue

Figure 7.3    Memory Assignment for Fixed Partitioning

## 2- Dynamic Partitioning

# 2-Dynamic Partitioning

- Partitions are of variable length and number.
- Process is allocated exactly as much memory as required.
- Eventually get holes in the memory. This is called **external fragmentation**.
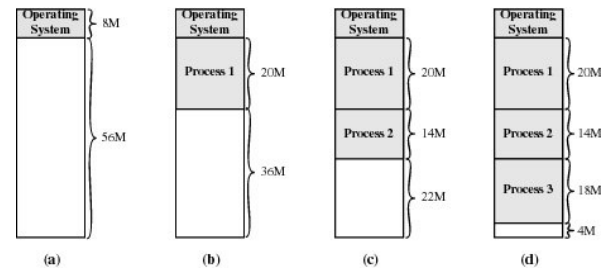- Must use **compaction** to shift processes so they are contiguous and all free memory is in one block.



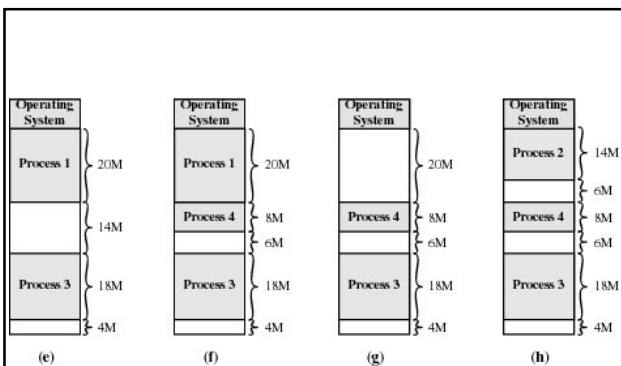Figure 7.4 The Effect of Dynamic Partitioning



Figure 7.4 The Effect of Dynamic Partitioning

# Memory Compaction

- To merge the memory blocks and reduce memory holes caused by external fragmentation.

- Time-consuming task....

## Dynamic Partitioning Placement Algorithms

- Operating system must decide which free block to allocate to a process.

- Four algorithms:
  **Best-Fit; First-Fit; Next-Fit; Worst-Fit**

## Dynamic Partitioning Placement Algorithm

- **Best-fit algorithm**
  - Chooses the block that is closest in size to the request.
  - Worst performance.
  - Since smallest block is found for process, the smallest amount of fragmentation is left memory compaction must be done more often.

## Dynamic Partitioning Placement Algorithm

- **Worst-fit algorithm**
  - Chooses the largest available block.
  - The fragmentation area can be used for another request.

## Dynamic Partitioning Placement Algorithm

- **First-fit algorithm**
  - Begins to scan memory from the beginning and chooses the next available block that is largest enough.
  - Fastest
  - May have many process loaded in the front end of memory that must be searched over when trying to find a free block.

Analysis of First-Fit Algorithm:

For every N allocated blocks, there are another 0.5N blocks that will be lost due to fragmentation -> 1/3 of memory will be unusable...

# Dynamic Partitioning Placement Algorithm

- **Next-fit**
  - Begins to scan memory from the location of the last placement and chooses the next available block that is largest enough.
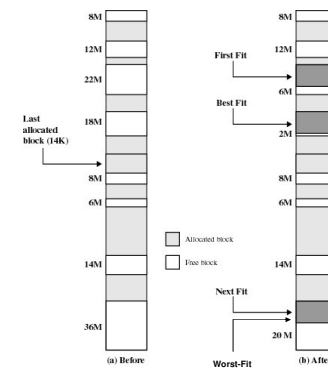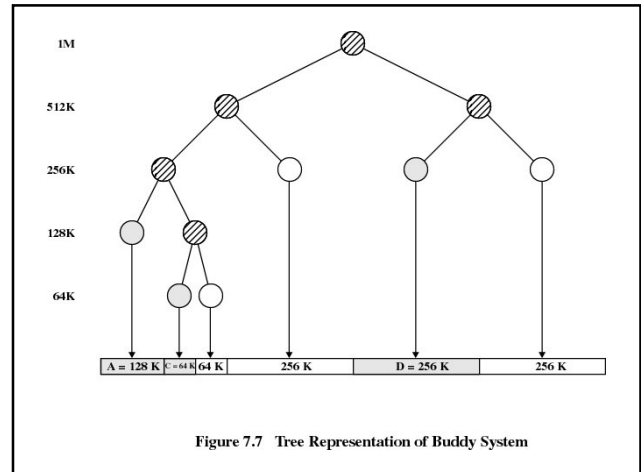


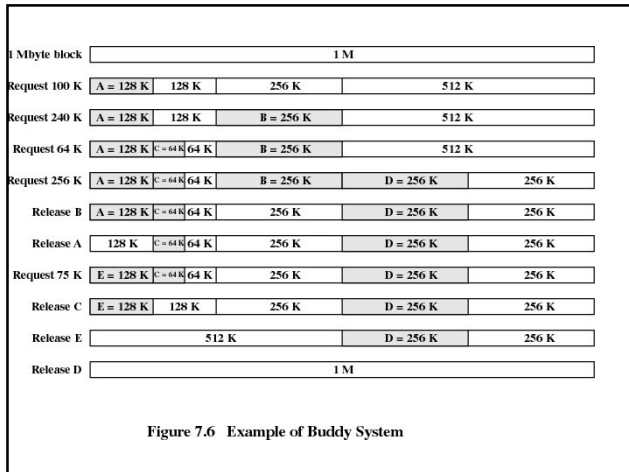Figure 7.5   Example Memory Configuration Before and After Allocation of 16 Mbyte Block

# 3-Buddy-System

# 3- Buddy System

- Entire space available is treated as a single block of $2^U$.
- If a request of size s such that $2^{U-1} < s <= 2^U$, entire block is allocated.
  - Otherwise block is split into two equal buddies.
  - Process continues until smallest block greater than or equal to s is generated.

Figure 7.6  Example of Buddy System



Figure 7.7  Tree Representation of Buddy System

# 4- Paging

# 4- Paging

- Partition memory into small equal-size chunks and divide each process into the same size chunks.
- The chunks of a process are called **pages** and chunks of memory are called **frames**.
- Operating system maintains a **page table** for each process
  - contains the frame location for each page in the process.
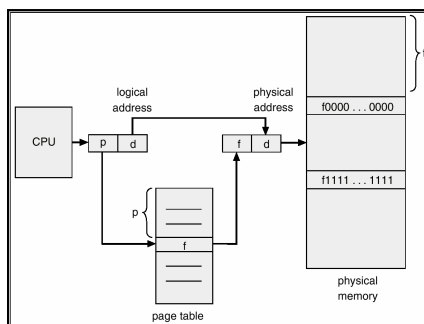  - memory address consist of a page number and offset within the page.

## Paging

- OS should 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 Scheme

- Address generated by CPU is divided into:
  - **Page number (p)** – used as an index into a *page table* which contains base address of each page in physical memory.

  - **Page offset (d)** – combined with base address to define the physical memory address that is sent to the memory unit.

## Paging: Address Translation



Page number (p) – used as an index into a *page table* which contains base address of each page in physical memory.
Page offset (d) – combined with base address to define the physical memory address that is sent to the memory unit.
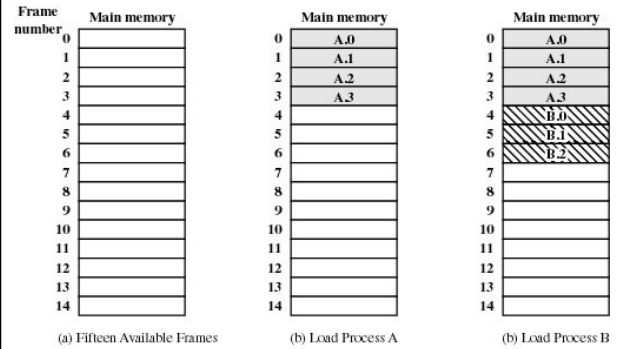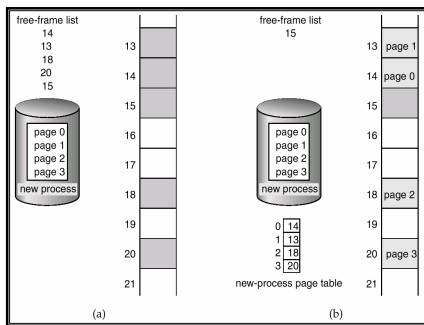
## Paging Example

# Free Frames

free-frame list
14
13
18
20
15

page 0
page 1
page 2
page 3
new process

13
14
15
16
17
18
19
20
21

(a)

free-frame list
15

page 0
page 1
page 2
page 3
new process

0 | 14
1 | 13
2 | 18
3 | 20
new-process page table

13 | page 1
14 | page 0
15
16
17
18 | page 2
19
20 | page 3
21

(b)

---

Frame number

Main memory (a) Fifteen Available Frames

Main memory (b) Load Process A

Main memory (b) Load Process B

**Figure 7.9  Assignment of Process Pages to Free Frames**

| Frame | (a) | (b) Load Process A | (b) Load Process B |
|---|---|---|---|
| 0 | | A.0 | A.0 |
| 1 | | A.1 | A.1 |
| 2 | | A.2 | A.2 |
| 3 | | A.3 | A.3 |
| 4 | | | B.0 |
| 5 | | | B.1 |
| 6 | | | B.2 |
| 7 | | | |
| 8 | | | |
| 9 | | | |
| 10 | | | |
| 11 | | | |
| 12 | | | |
| 13 | | | |
| 14 | | | |

---

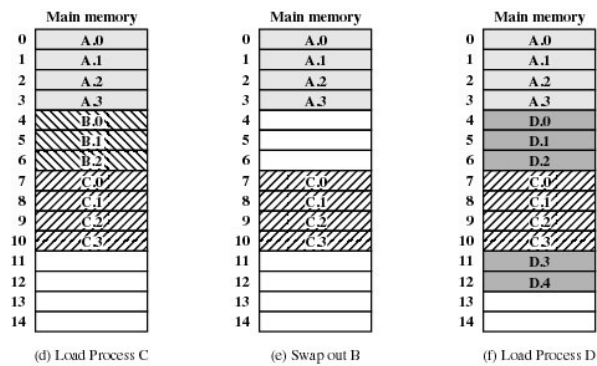| (d) Load Process C | (e) Swap out B | (f) Load Process D |
|---|---|---|
| 0 A.0 | 0 A.0 | 0 A.0 |
| 1 A.1 | 1 A.1 | 1 A.1 |
| 2 A.2 | 2 A.2 | 2 A.2 |
| 3 A.3 | 3 A.3 | 3 A.3 |
| 4 B.0 | 4 | 4 D.0 |
| 5 B.1 | 5 | 5 D.1 |
| 6 B.2 | 6 | 6 D.2 |
| 7 C.0 | 7 C.0 | 7 C.0 |
| 8 C.1 | 8 C.1 | 8 C.1 |
| 9 C.2 | 9 C.2 | 9 C.2 |
| 10 C.3 | 10 C.3 | 10 C.3 |
| 11 | 11 | 11 D.3 |
| 12 | 12 | 12 D.4 |
| 13 | 13 | 13 |
| 14 | 14 | 14 |

**Figure 7.9   Assignment of Process Pages to Free Frames**

---

# Page Tables for Example

Process A page table
0 | 0
1 | 1
2 | 2
3 | 3

Process B page table
0 | —
1 | —
2 | —

Process C page table
0 | 7
1 | 8
2 | 9
3 | 10

Process D page table
0 | 4
1 | 5
2 | 6
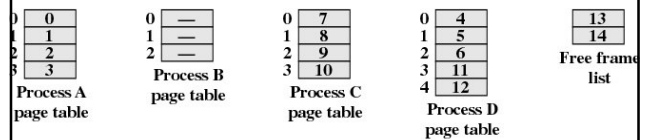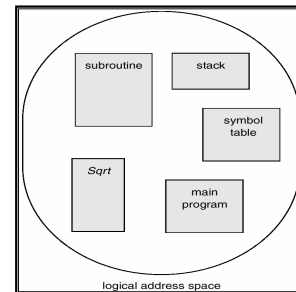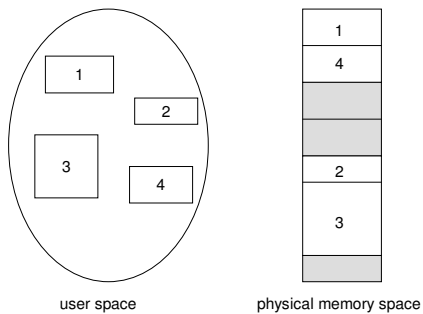3 | 11
4 | 12

Free frame list
13
14

**Figure 7.10  Data Structures for the Example of Figure 7.9 at Time Epoch (f)**

## 5- Segmentation

## User's View of a Program



subroutine
stack
symbol table
Sqrt
main program

logical address space

## Logical View of Segmentation



user space          physical memory space
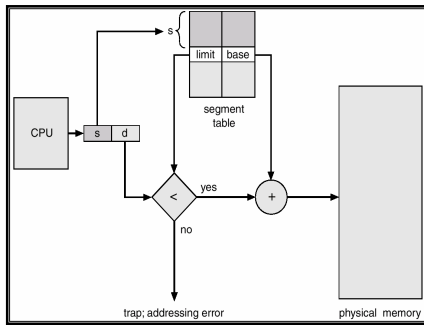
## 5- Segmentation

- All segments of all programs do not have to be of the same length.
- There is a maximum segment length.
- Addressing consist of two parts:
  - a segment number and an offset.
- Since segments are not equal, segmentation is similar to dynamic partitioning.
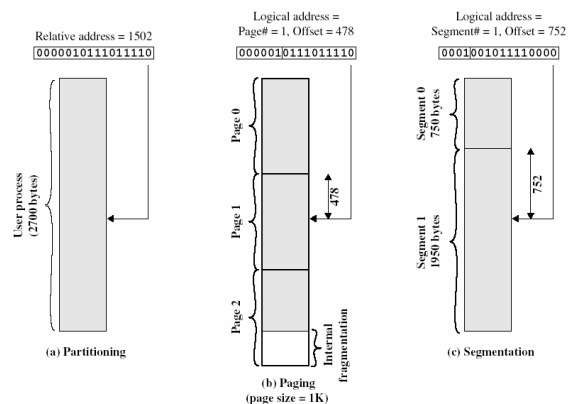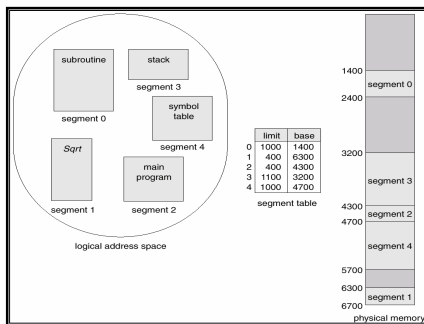
# Segmentation Hardware



# Segmentation Architecture
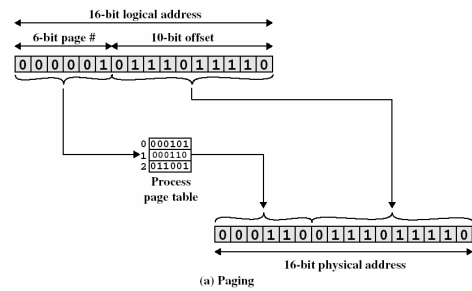
- Logical address consists of a two tuple:
- **<segment-number, offset>,**
- *Segment table* – maps two-dimensional physical addresses; each table entry has:
  - *base* – contains the starting physical address where the segments reside in memory.
  - *limit* – specifies the length of the segment.
- *Segment-table base register (STBR)* points to the segment table's location in memory.
- *Segment-table length register (STLR)* indicates number of segments used by a program;
  segment number $s$ is legal if $s <$ STLR.

# Example of Segmentation

# Logical-to-Physical Address Translation

## Logical-to-Physical Address Translation (Paging)

16-bit logical address

6-bit page #   10-bit offset

| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 |

0 | 000101
1 | 000110
2 | 011001
Process
page table

| 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 | 0 |

16-bit physical address

(a) Paging

## Logical-to-Physical Address Translation (Segmentation)

16-bit logical address

4-bit segment #   12-bit offset

| 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |

Length          Base
0 | 001011101110 | 0000010000000000
1 | 011110011110 | 0010000000100000
Process segment table

(+)

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

16-bit physical address

(b) Segmentation

16