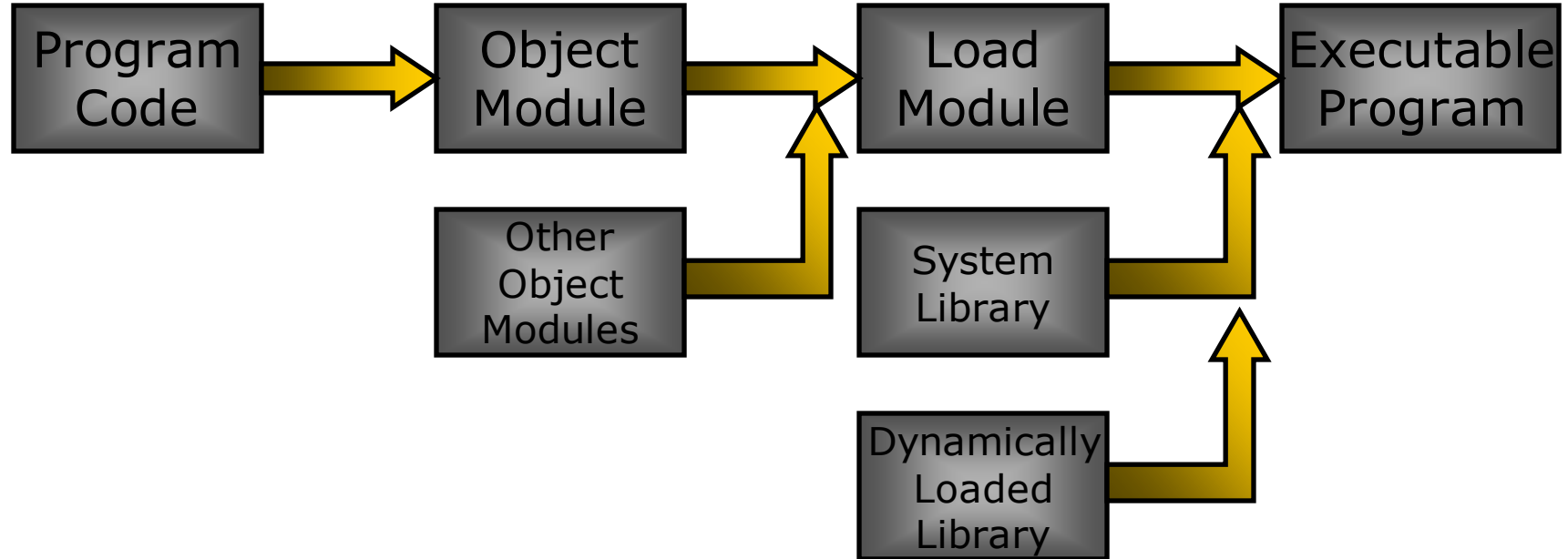


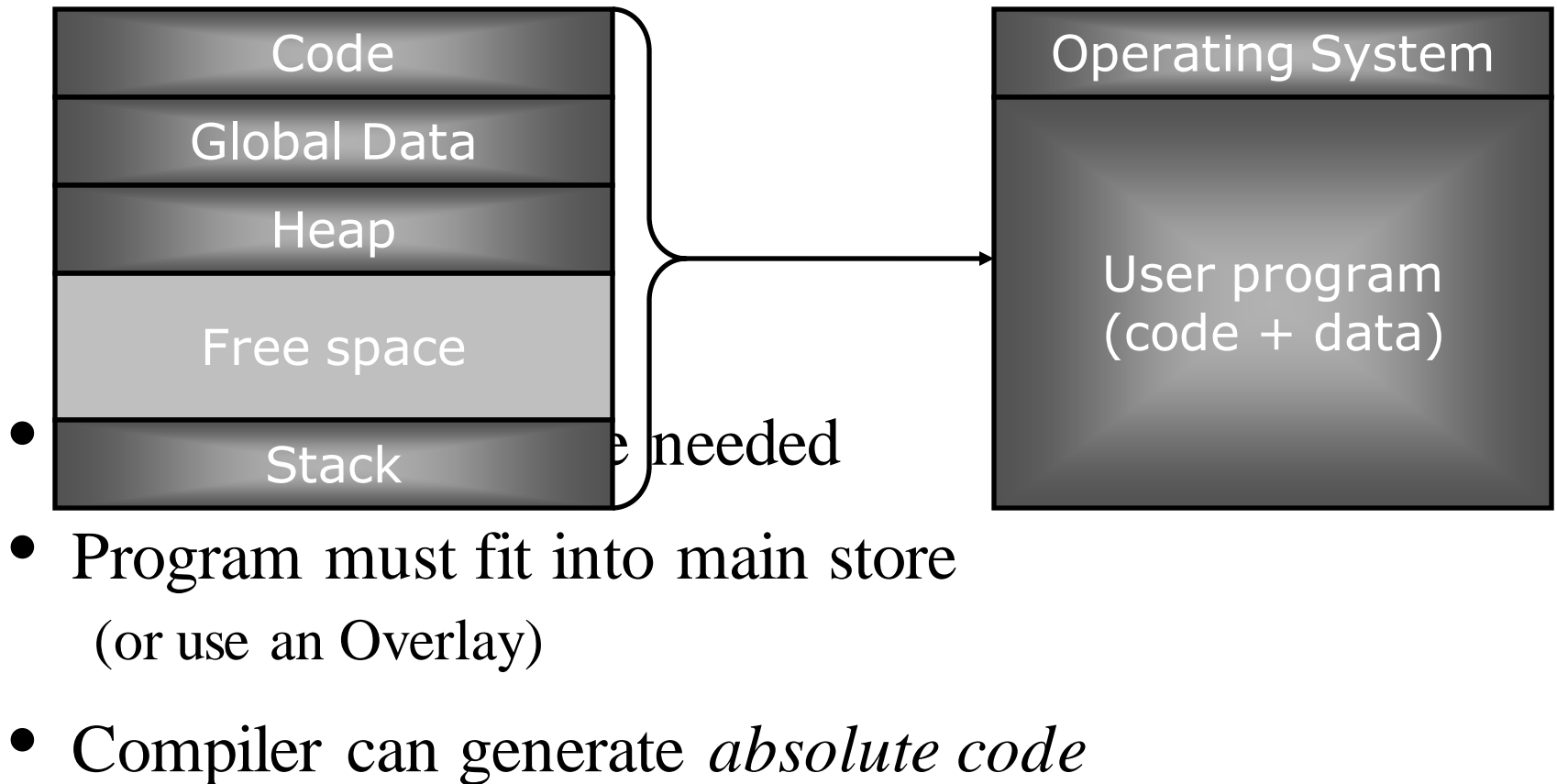
Lecture 9: Memory Management

- Memory management is about how the OS manages the allocation of memory to the various processes in a system
- Topics:
- Logical versus Physical Address Space
- Contiguous Allocation
- Paging
- Segmentation

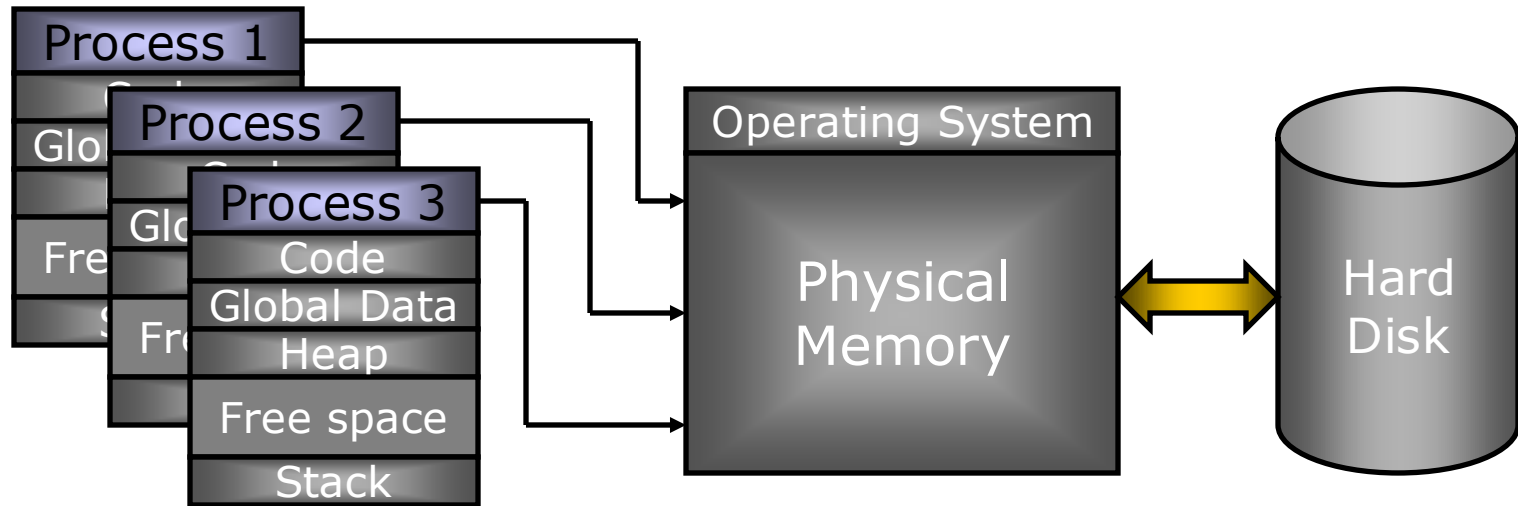
From Source Code to an Executing Program



Single Contiguous Store Allocation



Multiprocessing System



- Processes can reside anywhere in memory
 - Or even be swapped out to disk

Address binding

- Program must be brought into memory and converted into a process for it to be executed.
- More than one process is held in physical memory at the same time
- During execution processes make references to physical memory locations i.e. locations in chips
- BUT - Program code represents locations symbolically e.g. variable names
- Mapping of one type of address to another type of address known as **address binding**
- Address binding can happen at three different stages.

Binding of Instructions/Data to Memory

- 1. Compile time:** If the location of where the process is to be placed in memory is known when the program is compiled, then compiler can replace symbolic references to data in program code with actual addresses in memory where that data is to be held. You must recompile code if location of process in memory is to change.

2. Load time: If the location of where the process is to be placed in memory is NOT known when the program is compiled - compiler replaces symbolic references to data with **relocatable addresses**. These are addresses that are relative to address of beginning of program code I.e. first word of program code is given address 0 and then all other locations within program are given addresses as offset from start of program. The actual address of beginning of program code when it is placed in memory is added to all addresses when program is loaded into memory.

3. Execution time: the location of where the process is to be placed in memory is only known when the process is to be executed (run time binding) - as before compiler replaces symbolic references to data with **relocatable addresses**. However, process may be moved during its execution from one part of memory to another. Hardware (e.g., *base* and *limit registers*) is used to generate actual address to identify location in memory chip - e.g. relocatable address + address in some relocation register gives real address to be used. This is done for each memory reference.

Logical vs. Physical Address Space

- Execution time binding of addresses means that the address that is generated by the process when it is running on the CPU may be different to the address of the actual physical location that holds the instruction or data required - the address generated by the process running on CPU (known as **logical address** or **virtual address**) has to be translated into **physical address** of location in memory
- Logical and physical addresses are the same in compile-time and load-time address-binding, but may be different in run-time binding

Memory-Management Unit (MMU)

- MMU is a hardware device that maps virtual to physical address.
- In MMU every address generated by a user process is modified (e.g. by adding a relocation register value to address) before sending it to memory.
- Base and limit registers are an example of a simple MMU.

Memory allocation

- memory allocation is the process by which the OS decides what memory to allocate to a process for its code, data and system data
- There are 3 general schemes for allocating memory to processes:
 - contiguous allocation
 - paging
 - segmentation

Contiguous Allocation

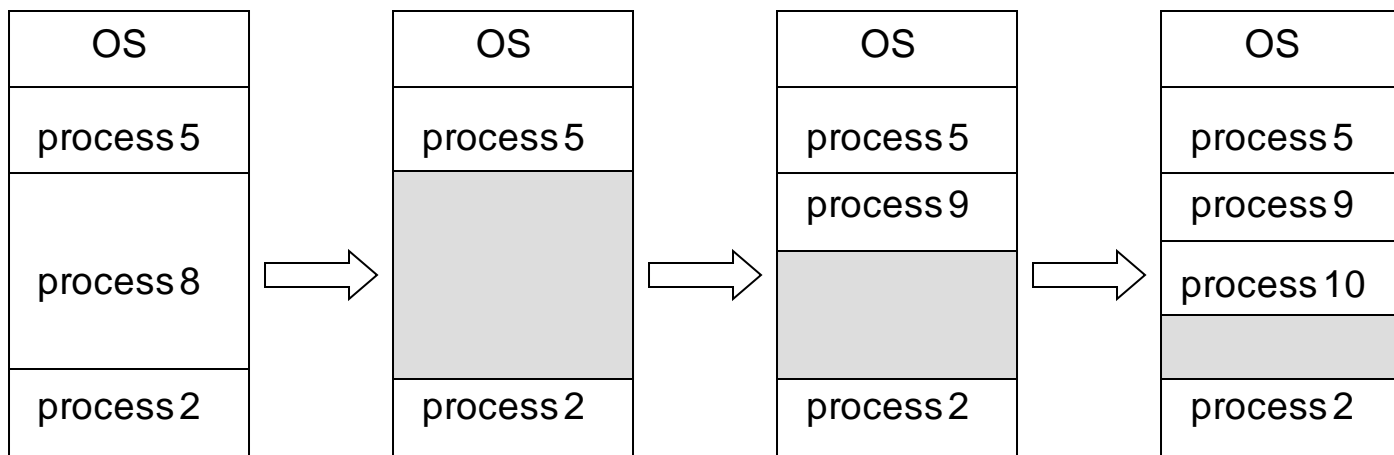
- The physical memory space of a process is held in one area of memory i.e. all locations are contiguous (next to each other)
- Main memory divided into two parts:
 - Resident operating system - held in low memory.
 - Memory for processes in high memory
- memory area for processes can be organised either
 - to hold the address space of a single process at a time (called single partition allocation)
 - to hold the address space of several processes at the same time (multiple partition allocation) - one partition for each process

Single-partition allocation

- Relocation-register scheme used to protect operating-system code and data.
- Base register (also sometimes called relocation register) contains value of physical address of start of process memory area, limit register contains range of logical addresses – each logical address must be less than the limit register.

Contiguous Allocation (Cont.)

- Multiple-partition allocation
 - **Hole** – this is a block of memory that is currently available (free) to be allocated to a process. Holes of various sizes 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 (holes)



Dynamic Storage-Allocation Problem

- How to satisfy a request for memory of size n from a list of free holes:
- **First-fit:** Allocate the *first* hole from list 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.

Fragmentation

- **External fragmentation** – enough total memory space exists for a process, but it is not all in one block i.e. not contiguous. The free memory is divided into small areas (many small holes) - thus process cannot be loaded into memory even though total free space is enough.
- The unusable memory space exists outside of the memory partitions allocated to the running processes, thus it is called external fragmentation.

- Reduce external fragmentation by compaction
 - Shuffle memory contents in memory to place all free memory together in one large block.
 - Compaction is possible *only* if relocation is dynamic, and is done at execution time.
- **Internal fragmentation** – memory allocated to a process may be slightly larger than amount of memory requested. This may be because if only the area of memory actually needed is allocated then the area leftover is too small to be of any use so all block is allocated. The unused area of memory is internal to a partition. Hence it is called internal fragmentation.

Paging

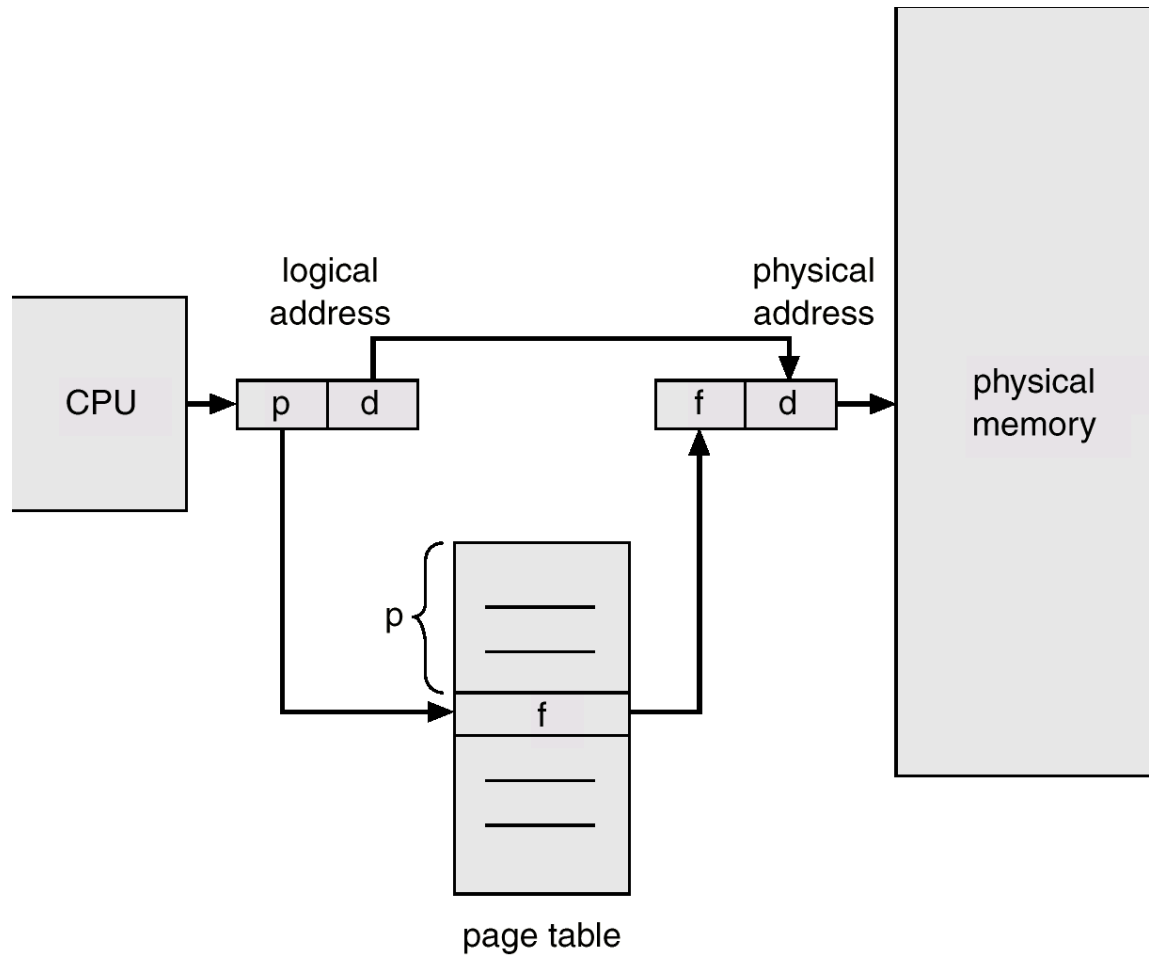
- In paging physical address space of a process can be noncontiguous i.e. addresses of consecutive logical locations in process may have non-consecutive physical addresses. The memory of the process may be divided up into a number of separate blocks in memory, the physical address is relative to the start address of the block of memory in which the required memory resides. As a result a process can be allocated physical memory wherever the latter is available.

- In paging scheme physical memory is divided up into into fixed-sized blocks called frames (normally between 1/2K and 8K in size).
- logical memory of a process is also divided into blocks of same size as the page frames. These blocks called pages.
- OS keeps track of all free frames.
- To run a program of size n pages, need to find n free frames and load program into them.
- Then set up a page table to translate logical to physical addresses. Separate page table for each process.

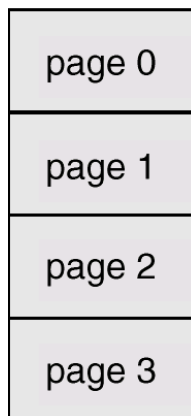
Address Translation Scheme

- 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 eliminates external fragmentation, but internal fragmentation can still exist.

Address Translation Architecture



Paging Example

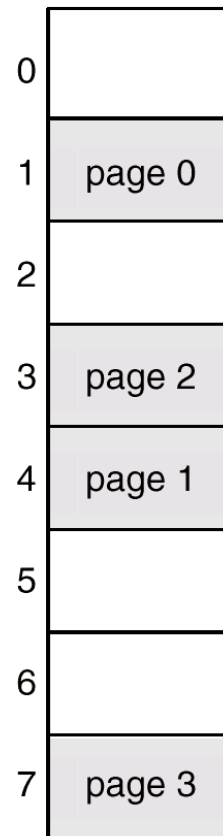


logical
memory

0	1
1	4
2	3
3	7

page table

frame
number



physical
memory

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 to get the page frame address from the page table and one to actually get the data/instruction.
- The two memory access problem can be solved by the use of a special fast-lookup hardware cache called *associative registers* or *translation look-aside buffers (TLBs)*

Associative Register

- Associative registers – allow parallel search

Page #	Frame #

- If page number is in associative register, get frame number out.
- Otherwise get frame number from page table in memory and add it to associative registers

Effective Access Time

- Associative Lookup = T_a time units
- Assume memory cycle time is 1 microsecond
- Hit ratio – percentage of times that a page number is found in the associative registers; ratio related to number of associative registers.
- If hit ratio = α , then Effective Access Time (EAT) can be found as follows:

if page no. in associative regs then use 1
microsecond to access memory for data + time
for associative lookup (this will happen α
fraction of the time)

if page no. not in associative regs then in addition
to time for associative lookup, use 2
microseconds to access page frame address and
data both in memory (this will happen $1 - \alpha$
fraction of the time)

$$\begin{aligned} \text{EAT} &= \alpha(1 + T_a) + (1 - \alpha)(2 + T_a) \\ &= 2 + T_a - \alpha \end{aligned}$$

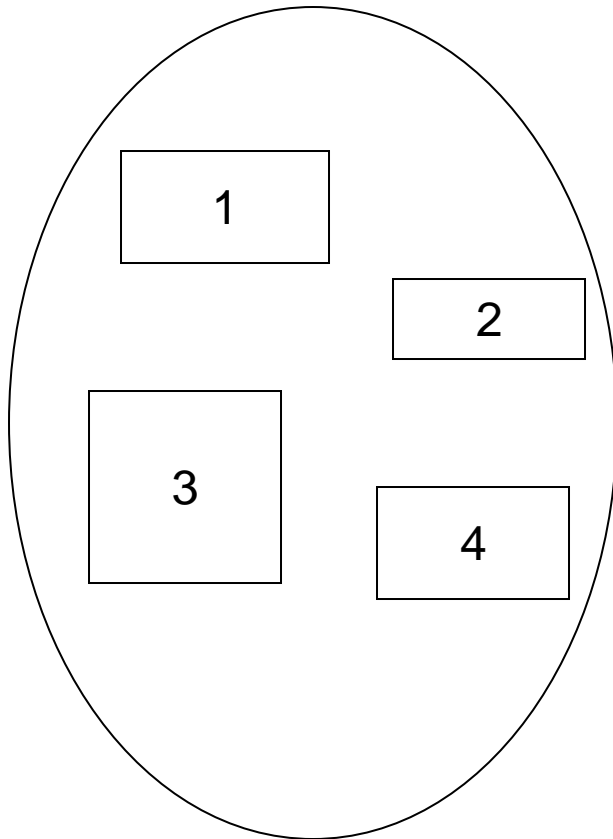
Memory Protection

- Memory protection is implemented by associating a protection bit with the page table entry for each frame - the bit specifies whether frame is read/write or read only
- Also *Valid-invalid* bit attached to each entry in the page table - bits loaded for each process as part of context switch for process:
 - “valid” indicates that the associated page is in the process’ logical address space, and is thus a legal page for the process to access.
 - “invalid” indicates that the page is not in the process’ logical address space and should thus generate an error

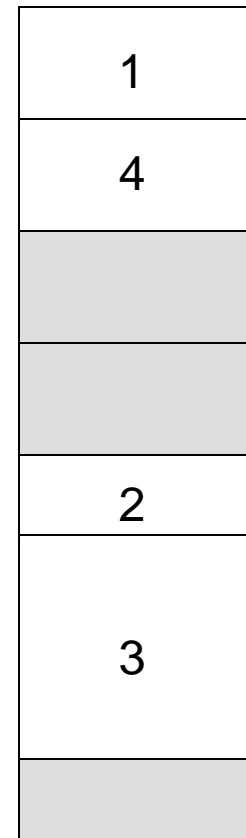
Segmentation

- Memory-management scheme that divides up processes' address space into a number of blocks called segments
- segments are similar to pages, but the difference is that pages all have the same size, whereas segments may be a number of different sizes
- In this scheme a program is treated as a collection of segments. A segment is a logical unit e.g. main part of program, a group of functions, the set of global data values, arrays, a stack, in OO languages may be individual objects

Logical View of Segmentation



user space



physical memory space

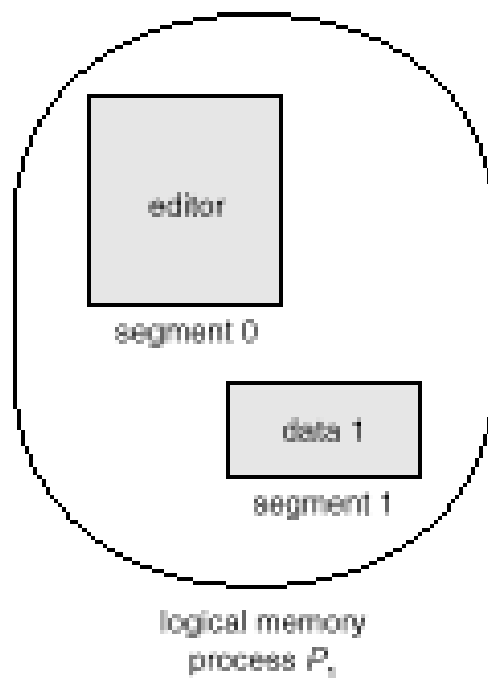
Segmentation Architecture

- Logical address consists of a pair: $\langle \text{segment-number}, \text{offset} \rangle$,
 - Segment number is an index into the *segment table* - with one segment table for each process
 - offset is the value to be added to start address of segment to give real physical address
- each entry in segment table has:
 - base value – contains the starting physical address where the segment resides in memory.
 - *limit* – specifies the length of the segment.
- Offset address is legal if $\text{offset} < \text{limit value}$

- *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 $s < \text{STLR}$.

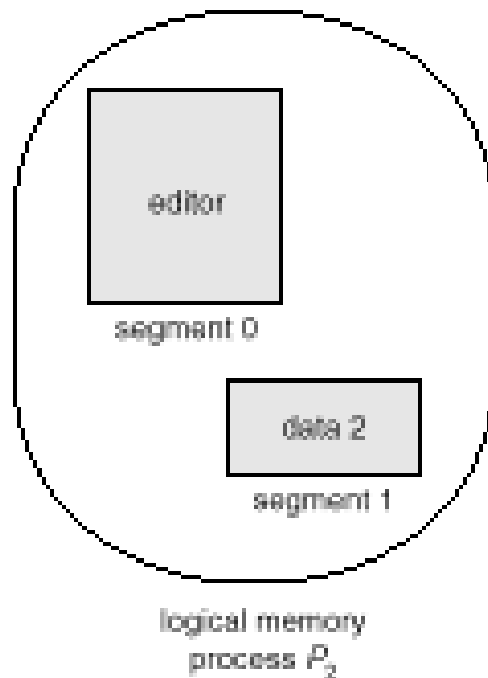
Segmentation Architecture (Cont.)

- Protection. Each entry in segment table contains:
 - validation bit - which if = 0 \Rightarrow illegal segment
 - specification of read/write/execute privileges
- Protection bits associated with segments
- Since segments vary in length, memory allocation is a dynamic storage-allocation problem.
- A segmentation example is shown in the following diagram



	limit	base
0	25286	43062
1	4425	68348

segment table
process P_1



	limit	base
0	25286	43062
1	8850	90003

segment table
process P_2

