# Lesson 18

# **Objectives**

- Memory Management
- Address Binding
- Logical Vs Physical Address Space
- MMU
- Dynamic Loading
- Dynamic Linking
- Overlays
- Swapping

# **Memory Management**

Memory is the second important resource in computer system. So a good operating system should have an optimum memory management plan to increase the throughput of the system since even performance of CPU depends upon performance of main memory.

- Program must be brought into memory and placed within a process for it to be run.
- Input queue collection of processes on the disk that are waiting to be brought into memory to run the program.
- User programs go through several steps before being run.

## **Address Binding**

Address binding of instructions and data to memory addresses can happen at three different stages.

## 1. Compile time:

If memory location known a priori, absolute code can be generated; must recompile code if starting location changes.

#### 2. Load time:

Must generate re-locatable code if memory location is not known at compile time.

#### 3. Execution time:

Binding delayed until run time if the process can be moved during its execution from one memory segment to another. Need hardware support for address maps (e.g., base and limit registers).

These different times are depicted in the figure below. The moment program is being assembled and compiled that time is referred as Compile Time. Linking and loading time is referred Loading Time, and the moment program is being executed is called Execution Time.

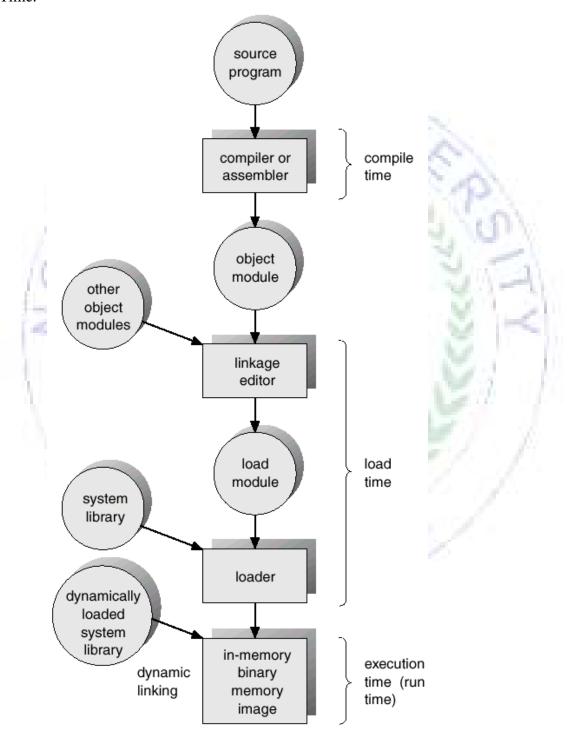


Figure: Different memory allocation times

## **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.

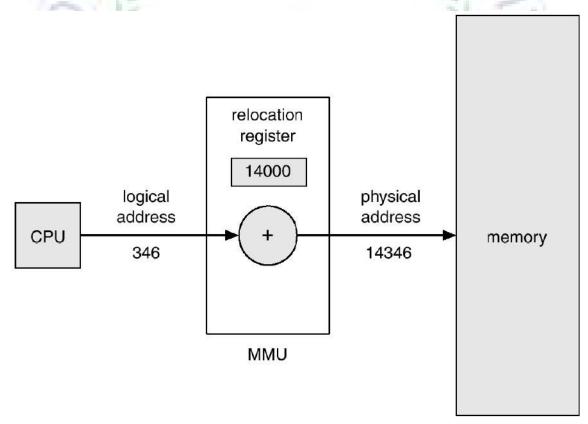
- Logical address generated by the CPU; also referred to as virtual address.
- Physical address address seen by the memory unit

Logical addresses are converted to physical address, that process is called address translation or address mapping. That is done by Memory Management Unit (MMU).

# **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.



**Mapping Process** 

# **Dynamic Loading**

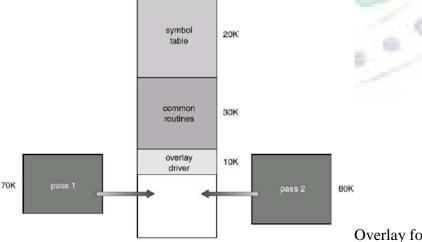
- Routine is not loaded until it is called
- Better memory-space utilization; unused routine is never loaded.
- Useful when large amounts of code are needed to handle infrequently occurring cases.
- No special support from the operating system is required implemented through program design.

# **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 needed to check if routine is in processes' memory address.
- Dynamic linking is particularly useful for libraries.

# Overlays

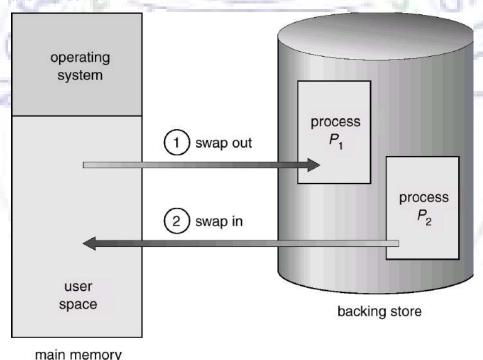
- Keep in memory only those instructions and data that are needed at any given time.
- Needed when process is larger than amount of memory allocated to it.
- Implemented by user, no special support needed from operating system,
  programming design of overlay
- structure is complex



Overlay for a two pass assembler

# **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.



# **Memory Management**

Main Memory refers to a physical memory that is the internal memory to the computer. The word main is used to distinguish it from external mass storage devices such as disk drives. Main memory is also known as RAM. The computer is able to change only data

that is in main memory. Therefore, every program we execute and every file we access must be copied from a storage device into main memory.

All the programs are loaded in the main memory for execution. Sometimes complete program is loaded into the memory, but some times a certain part or routine of the program is loaded into the main memory only when it is called by the program, this mechanism is called **Dynamic Loading**, this enhance the performance.

Also, at times one program is dependent on some other program. In such a case, rather than loading all the dependent programs, CPU links the dependent programs to the main executing program when it's required. This mechanism is known as **Dynamic Linking**.

## **Swapping**

A process needs to be in memory for execution. But sometimes there is not enough main memory to hold all the currently active processes in a timesharing system. So, excess process are kept on disk and brought in to run dynamically. Swapping is the process of bringing in each process in main memory, running it for a while and then putting it back to the disk.

# **Contiguous Memory Allocation**

In contiguous memory allocation each process is contained in a single contiguous block of memory. Memory is divided into several fixed size partitions. Each partition contains exactly one process. When a partition is free, a process is selected from the input queue and loaded into it. The free blocks of memory are known as *holes*. The set of holes is searched to determine which hole is best to allocate.

## **Memory Protection**

Memory protection is a phenomenon by which we control memory access rights on a computer. The main aim of it is to prevent a process from accessing memory that has not been allocated to it. Hence prevents a bug within a process from affecting other processes, or the operating system itself, and instead results in a segmentation fault or storage violation exception being sent to the disturbing process, generally killing of process.

## **Memory Allocation**

Memory allocation is a process by which computer programs are assigned memory or space. It is of three types:

## 1. First Fit

The first hole that is big enough is allocated to program.

#### 2. **Best Fit**

The smallest hole that is big enough is allocated to program.

#### 3. Worst Fit

The largest hole that is big enough is allocated to program.

## **Fragmentation**

Fragmentation occurs in a dynamic memory allocation system when most of the free blocks are too small to satisfy any request. It is generally termed as inability to use the available memory.

In such situation processes are loaded and removed from the memory. As a result of this, free holes exists to satisfy a request but is non-contiguous i.e. the memory is fragmented into large no. Of small holes. This phenomenon is known as **External Fragmentation.** 

Also, at times the physical memory is broken into fixed size blocks and memory is allocated in unit of block sizes. The memory allocated to a space may be slightly larger than the requested memory. The difference between allocated and required memory is known as **Internal fragmentation** i.e. the memory that is internal to a partition but is of no use.

#### **Paging**

A solution to fragmentation problem is Paging. Paging is a memory management mechanism that allows the physical address space of a process to be non-contagious. Here physical memory is divided into blocks of equal size called **Pages**. The pages belonging to a certain process are loaded into available memory frames.

## Page Table

A Page Table is the data structure used by a virtual memory system in a computer operating system to store the mapping between *virtual address* and *physical addresses*.

Virtual address is also known as Logical address and is generated by the CPU. While Physical address is the address that actually exists on memory.

## Segmentation

Segmentation is another memory management scheme that supports the user-view of memory. Segmentation allows breaking of the virtual address space of a single process into segments that may be placed in non-contiguous areas of physical memory.

# **Segmentation with Paging**

Both paging and segmentation have their advantages and disadvantages, it is better to combine these two schemes to improve on each. The combined scheme is known as 'Page the Elements'. Each segment in this scheme is divided into pages and each segment is maintained in a page table. So the logical address is divided into following 3 parts:

- Segment numbers(S)
- Page number (P)
- The displacement or offset number (D)

