

Operating System – Unit 3

Topic: Memory Management

Memory: Memory in the context of computers refers to the physical devices or components that store data and instructions for processing by the CPU (Central Processing Unit). It's like the workspace where a computer keeps information temporarily, allowing it to perform tasks and execute programs.

Types of Memory:

1. Primary Memory (RAM - Random Access Memory):

- **Function:** Stores data and instructions that the CPU is currently using.
- **Characteristics:** Fast, volatile (loses data when powered off), and directly accessed by the CPU.

2. Secondary Memory (Storage Devices - HDD, SSD, etc.):

- **Function:** Holds data even when the computer is turned off.
- **Characteristics:** Slower than primary memory, non-volatile (retains data when powered off), and used for long-term storage.

Purpose:

- Memory allows the computer to store and retrieve data quickly, facilitating multitasking, running programs, and storing both temporary and permanent data crucial for the computer's operation.

Discuss various Memory Management techniques or Schemes in detail.

Memory Management Schemes are techniques used by operating systems to efficiently handle and organize memory resources. These schemes vary in how they allocate, utilize, and manage memory for running programs.

Here are the following Memory Management Schemes or techniques:

1. Contiguous, Real Memory Management System:

- **Contiguous Allocation:** In this method, memory is divided into fixed-sized partitions. Each partition accommodates one process entirely. Programs are loaded into contiguous memory blocks. For instance, if a program needs 100 KB and a partition is 200 KB, the whole partition will be allocated.
- **Advantages:** Simplicity in management and memory access, minimal overhead in allocation and deallocation.
- **Disadvantages:** Fragmentation issues can arise, particularly external fragmentation where free memory exists but scattered in small chunks, causing wastage.

2. Non-contiguous, Real Memory Management System:

- **Dynamic Partitioning:** Memory is allocated dynamically based on the actual size of programs. Allocation happens in non-contiguous blocks, meaning different parts of a program can occupy separate memory locations.
- **Advantages:** Better space utilization and flexibility in accommodating varying program sizes. Reduces external fragmentation by allocating memory based on actual needs.
- **Disadvantages:** More complex memory management, requires dynamic allocation and deallocation, leading to increased overhead.

3. Non-Contiguous, Virtual Memory Management:

- **Virtual Address Space:** Programs access an extended logical address space beyond physical memory limits. The system manages this by using secondary storage (like a hard drive) as an extension of RAM.
- **Paging or Segmentation:** Memory is divided into fixed-size pages or logical segments. Pages or segments are swapped between RAM and secondary storage based on the program's demand.
- **Advantages:** Efficiently executes larger programs by using secondary storage as an extension of physical memory.

- **Disadvantages:** Slower access compared to physical memory due to disk access times. Requires complex mapping mechanisms to translate virtual addresses to physical ones.

Memory Address: A Memory Address is a unique identifier assigned to each location in a computer's memory. It acts like a street address, pinpointing a specific location where data or instructions are stored. Each byte of memory has its own unique address, allowing the computer to access and retrieve information quickly and accurately.

1. Absolute Address:

- **Simple Explanation:** An Absolute Address is like a unique street address for data in a computer's memory.
- **More Details:** Just as your home has a specific address that uniquely identifies its location, an Absolute Address points directly to a particular spot in the computer's memory. It's a fixed and permanent reference that doesn't change, making it easy for the computer to find data quickly. However, if a program needs to move to a different memory location, using Absolute Addresses can be challenging because they're tied to specific locations.

2. Relative Address:

- **Simple Explanation:** A Relative Address is like giving directions from a landmark rather than using an exact address.
- **More Details:** Instead of providing an exact location, a Relative Address tells you how far or where something is from a specific starting point or reference. It's like saying, "The store is three blocks away from the library." This type of address doesn't rely on fixed locations but rather describes distances or offsets from a reference point. It's helpful when programs need to move around in memory because Relative Addresses are adaptable and don't depend on fixed memory locations.

Address Binding: Address binding is the process of connecting or associating a symbolic identifier (like a variable name in a program) with a specific memory location or address in a computer's memory. This association can happen at different stages in the program's lifecycle, depending on the memory management scheme used by the operating system.

There are mainly three types of address binding:

1. **Compile Time (Static Binding):** Address binding occurs when the program is compiled. The addresses for variables and instructions are fixed and determined before the program starts running. This binding type is rigid, and the program cannot change memory locations once compiled. It's efficient but lacks flexibility.
2. **Load Time (Static or Partial Binding):** Address binding happens when the program is loaded into memory for execution. Memory addresses for variables and instructions are assigned when the program is loaded, but some parts may remain flexible, allowing the program to adapt to different memory locations. However, the core addresses are fixed once the program starts running.
3. **Execution Time (Dynamic or Runtime Binding):** Address binding occurs while the program is running. The addresses for variables and instructions are determined dynamically during runtime. This method offers the most flexibility as memory locations can change dynamically, but it requires additional runtime processing, making it slower compared to other binding types.

Explain the three phases of program execution ?

1. **Compilation Phase:** During this phase, the program's source code, written in high-level programming languages like C, C++, Java, etc., is translated into machine code or object code that the computer can understand.
 - **Process:** A compiler or an interpreter converts the human-readable source code into a format that the computer's CPU can execute. The output is often in the form of object files or executable files.

2. **Loading Phase:** Once the program has been compiled into object code or an executable file, the Loading Phase comes into play.
 - **Process:** In this phase, the operating system loads the compiled program from the storage (like a hard drive or solid-state drive) into the computer's memory (RAM) for execution. It assigns memory addresses to the program's instructions and data so that the CPU can access and execute them.
3. **Execution Phase:** This is the phase where the actual execution of the program takes place.
 - **Process:** The CPU begins executing the loaded program instructions one by one in the order specified. Data processing, calculations, and any other operations defined in the program are carried out. The program continues its execution until it completes its tasks, encounters an error, or is manually stopped.

Explain Memory Mapping.

Memory mapping involves converting logical addresses, generated by a program during its execution, into corresponding physical addresses in the computer's memory. This translation or mapping process is handled by the operating system in coordination with the hardware components, primarily to manage memory efficiently and facilitate the execution of programs.

Here's a simplified explanation:

1. Physical Address: A Physical Address is like the actual street address of a house.

It's the real location in the computer's memory where data is stored, just like a house's actual location on a street. Computers use these addresses to directly access and retrieve data from memory chips, like finding a specific house on a street.

2. Logical Address: Logical Address is like a virtual address, similar to using a GPS location rather than the actual street address.

When a program runs, it uses logical addresses, thinking it's accessing specific memory locations. However, these addresses aren't the real physical locations; they're like GPS coordinates that the computer translates into the actual physical addresses to find the data.

3 Mapping Process: The operating system manages the translation of logical addresses into physical addresses. It does this by using a Memory Management Unit (MMU), a hardware component that maps or translates logical addresses into their corresponding physical addresses.

4. Translation Mechanism: The MMU translates each logical address generated by the program into its actual physical memory location. It uses various techniques, like page tables (in the case of virtual memory), to perform this translation efficiently.

5. Purpose: The primary goal of memory mapping is to enable programs to work with logical addresses while efficiently mapping them to the physical memory. This mapping ensures that programs can access the correct memory locations without needing to be aware of the underlying physical memory organization.

Swapping: Swapping is a memory management technique used by operating systems to temporarily move parts of a program from the main memory (RAM) to secondary storage, such as a hard drive, when the RAM is insufficient to hold all the currently executing processes.

In a multiprogramming environment, where multiple programs or processes run concurrently, swapping plays a crucial role in managing the limited memory resources efficiently.

As multiple programs run simultaneously, their memory demands fluctuate. Some programs may require more memory at specific times, while others may be temporarily idle or use fewer resources.

Swapping enables the operating system to temporarily move less active or idle processes from the main memory (RAM) to the secondary storage (like a hard

drive) when the available RAM is insufficient to accommodate all active processes.

By swapping out less critical or inactive processes to the secondary storage, the operating system frees up space in the RAM for other more active processes to run. This helps in maximizing the utilization of the available memory resources.

Memory Allocation: Memory allocation refers to the process of assigning and managing memory resources to computer programs or processes during their execution. It involves dividing the available memory space among various tasks efficiently, ensuring that each program gets the necessary memory for execution while preventing conflicts or overlaps.

1. Static Memory Allocation:

- **Explanation:** Static memory allocation occurs during the compilation or linking phase of a program. Memory for variables or data structures is allocated and fixed at compile time and remains unchanged throughout the program's execution. It assigns memory addresses to variables, defining their size and location in memory before the program runs. However, this method lacks flexibility as the allocated memory cannot be adjusted during runtime.

2. Dynamic Memory Allocation:

- **Explanation:** Dynamic memory allocation happens during the execution of a program, allowing memory allocation and deallocation at runtime. Programs request memory as needed using functions like **malloc()** and **free()** in languages like C or C++. It offers flexibility as memory can be allocated or released dynamically based on program requirements, enabling efficient memory utilization and adaptability.

Partitions: Partitions in the context of computer memory refer to the division or segmentation of the available memory space into separate sections or blocks. These sections are used to allocate memory to different programs or processes

running concurrently in the system. There are primarily two types of memory partitioning:

1. Fixed Partitioning:

- **Explanation:** Fixed Partitioning involves dividing the available memory into fixed-sized partitions or segments during system initialization.
- **Characteristics:** Each partition is of a predetermined size, creating several fixed compartments in memory.
- **Usage:** When a process is loaded into memory, it is allocated to a partition that can accommodate its size. Only one process can occupy a partition at a time. If a process's size is smaller than the partition, there might be internal fragmentation (unused space within a partition).

1.1 Single partition allocation: Single partition allocation, often referred to as monoprogramming, represents a memory management scheme where the entire available memory is allocated to a single program or process. In this scenario, the operating system loads and executes only one program at a time, utilizing the entire memory space exclusively for that program's execution.

Here's a more detailed explanation:

Entire Memory for Single Program:

Explanation: In monoprogramming, the entire memory, including RAM, is dedicated to a single program without any division or partitioning.

Characteristics: The entire memory space serves the needs of the running program, allowing it to execute without sharing resources with other processes.

Utilization of Entire Resources:

Usage: When a program is loaded, it occupies the complete memory space available in the system, utilizing all the memory resources exclusively for its execution.

No Sharing: There are no other concurrent processes or programs running alongside the currently executed program.

Simplicity but Limitations:

Advantages: Monoprogramming offers simplicity in managing memory, as there's only one program to handle. This simplicity results in straightforward memory allocation and management.

Disadvantages: However, it severely limits the multitasking capability of the system since only one program can execute at a time. Other programs or tasks have to wait until the running program finishes its execution.

1.2 Multiple Fixed Partitions Allocation: Multiple Fixed Partitions Allocation is a memory management scheme where the available memory is divided into several fixed-size partitions or blocks to accommodate multiple programs or processes simultaneously. Each partition is of a predetermined size and can hold a single process at a time.

2. Dynamic Partitioning or Multiple Variable Partition Allocation:

- **Explanation:** Dynamic Partitioning allows the division of memory into variable-sized partitions to better fit processes of different sizes.
- **Characteristics:** Partitions are created and resized dynamically based on the size of the processes being executed.
- **Usage:** When a process is loaded, it is allocated memory based on its size, and the system searches for a suitable partition that can accommodate it. After a process finishes execution, its partition becomes available for new processes. Dynamic partitioning reduces internal fragmentation but may result in external fragmentation (free memory scattered in small blocks).

Explain Internal Fragmentation ?

Internal Fragmentation occurs in memory management when allocated memory space is not fully utilized, resulting in wastage of available memory due to the inefficient allocation of memory blocks to processes.

Here's a simplified explanation:

1. Memory Block Allocation:

- **Fixed Partitioning or Block Sizes:** In scenarios like fixed partitioning, memory is divided into fixed-size blocks or partitions.
- **Allocating Memory Blocks:** When a process is allocated memory, it might not perfectly fit the entire block allocated to it.

2. Unused Memory Within Blocks:

- **Unused Space:** If a process doesn't fully utilize the entire memory block allocated to it, the remaining space within that block remains unused.
- **Wasted Memory:** This leftover or unused space within the allocated block is referred to as internal fragmentation.

3. Causes and Impact:

Internal fragmentation occurs due to the mismatch between the memory block size and the actual memory requirement of the process.

4. Management and Mitigation:

- **Optimizing Allocation:** Techniques like dynamic partitioning or memory compaction aim to minimize internal fragmentation.
- **Memory Utilization:** Reducing internal fragmentation enhances overall memory utilization and system efficiency.

Explain External Fragmentation ?

External Fragmentation occurs in memory management when there's unallocated or free memory scattered in small blocks between allocated memory segments or partitions, making it challenging to utilize these small free spaces to accommodate larger processes, despite having enough total free memory.

Here's an easy explanation:

1. Free Memory Gaps:

- **Unused Memory Spaces:** As processes are loaded and removed from memory, small gaps or free spaces emerge between allocated memory blocks or partitions.

- **Scattered Free Blocks:** These free memory spaces might be too small individually to satisfy the memory requirements of incoming processes.

2. Impact on Memory Utilization:

- **Reduced Efficiency:** External fragmentation limits memory utilization efficiency, as free memory remains unused due to the inability to merge these small free blocks into a larger contiguous block.
- **Unavailable for Allocation:** Although the total amount of free memory might be sufficient, these scattered small free spaces cannot accommodate larger processes due to their fragmented nature.

3. Management Strategies:

- **Dynamic Memory Allocation:** Techniques like dynamic partitioning or virtual memory management aim to reduce external fragmentation by more flexible allocation strategies.
- **Unavailable for Allocation:** Although the total amount of free memory might be sufficient, these scattered small free spaces cannot accommodate larger processes due to their fragmented nature.

Explain Compaction ?

Compaction is like organizing a room with scattered toys. In the world of computers, it's a memory management technique used by the operating system to tidy up memory space. Imagine rearranging scattered Lego blocks to create larger, continuous spaces. Similarly, compaction gathers and merges smaller unused memory areas, creating bigger blocks of free space. This process helps prevent wasted memory caused by fragmentation, where small gaps of free memory exist between allocated spaces. By organizing memory more efficiently, compaction allows larger programs to fit into the available memory, improving the overall usage of the system's resources.

Topic: Virtual Memory and Paging

Explain Virtual Memory ?

Virtual Memory is a memory management technique used by operating systems to expand the available memory beyond the physical RAM (Random Access Memory) by using a portion of the hard disk space as an extension of the main memory. It allows the system to run programs larger than the physical memory size and provides a way to efficiently manage memory space.

Here's a simple breakdown:

1. **Memory Extension:** Virtual Memory acts as an extension of the physical RAM. It makes the system believe it has more memory than it physically possesses by temporarily transferring data from RAM to the hard disk.
2. **Address Space Division:** The memory space is divided into smaller parts called "pages." These pages are mapped between physical memory (RAM) and the hard drive.
3. **Page Swapping:** Virtual Memory swaps pages between RAM and the disk based on demand, ensuring active processes have access to required data while keeping the less used data on the disk.
4. **Advantages:**
 - **Larger Program Execution:** Allows running larger programs or multiple programs simultaneously, even when physical memory is limited.
 - **Enhanced System Stability:** Helps prevent system crashes due to lack of memory by providing a larger memory space.
5. **Disadvantages:**
 - **Performance Impact:** Accessing data from the hard disk is slower compared to RAM, which can affect system performance when frequent swapping occurs.
 - **Storage Limitation:** The size of virtual memory is limited by the available space on the hard disk.

Explain Paging ?

Paging is a memory management scheme used by operating systems to efficiently manage and allocate memory in a more flexible and adaptable manner compared to fixed partitioning. It involves dividing the physical memory into fixed-size blocks called "pages" and the logical memory into blocks of the same size known as "frames." These pages and frames allow for more efficient use of memory space and facilitate virtual memory management.

Here's a breakdown of paging:

1. Memory Division into Pages and Frames:

- **Fixed-Size Blocks:** Physical memory (RAM) is divided into fixed-size blocks known as "frames."
- **Corresponding Logical Blocks:** The logical memory space visible to programs is divided into blocks of the same size called "pages."

2. Page Table for Address Translation:

- **Address Translation:** Each program's logical address space is divided into pages with a unique page number.
- **Page Table:** A page table is maintained by the operating system, which maps logical page numbers to corresponding physical frame numbers, enabling efficient address translation from logical to physical memory addresses.

3. Flexible Memory Allocation:

- **Dynamic Allocation:** Pages of a process do not need to be contiguous in physical memory.
- **Non-contiguous Pages:** Pages of a process can reside in different non-contiguous frames, allowing for more flexible memory allocation.

4. Advantages of Paging:

- **Reduced Fragmentation:** Paging reduces external fragmentation as it doesn't require contiguous memory allocation for processes.

- **Efficient Use of Memory:** It allows for more efficient memory management by breaking memory into smaller, fixed-size units.

Advantages of Paging:

1. **Reduced Fragmentation:** Paging helps reduce external fragmentation compared to other memory allocation techniques like fixed partitioning. It allocates memory in smaller, fixed-sized units (pages), reducing the likelihood of unusable small memory gaps.
2. **Flexible Memory Allocation:** It offers flexibility in memory allocation as pages of a process don't need to be contiguous in physical memory. Processes can occupy non-contiguous frames, enabling more efficient memory utilization.
3. **Simplified Address Translation:** Paging simplifies address translation by using a page table. This translation mechanism efficiently converts logical addresses used by the program to physical addresses in memory.
4. **Support for Virtual Memory:** Paging forms the basis for virtual memory systems, allowing the system to manage more memory than physically available by swapping pages in and out of secondary storage (like a hard drive).

Disadvantages of Paging:

1. **Internal Fragmentation:** Paging can lead to internal fragmentation within pages if a process doesn't fully occupy an entire page. This unused space within a page can result in wasted memory.
2. **Overhead:** Maintaining and managing the page table for address translation incurs additional overhead on the system's resources, especially as the size of the address space increases.
3. **Complexity in Page Replacement:** Implementing efficient page replacement algorithms, which decide which pages to swap in and out of memory, can be complex and may impact system performance.

Explain Demand Paging ?

Demand Paging is a memory management strategy used by operating systems to optimize memory usage by loading only necessary parts of a program into memory when needed.

It works like borrowing a book from a library: instead of bringing the entire bookshelf home, you fetch and read only the specific pages you require. Similarly, with demand paging, when a program starts, only a small portion, or even just a single page, is initially loaded into memory. As the program runs and requires more data, the operating system fetches additional pages from the disk into memory on an as-needed basis. This "on-demand" fetching minimizes the initial loading time and allows the system to efficiently utilize memory resources by bringing in only the required parts of the program.

However, if a page that's needed is not in memory, it leads to a page fault, triggering the OS to fetch that specific page from secondary storage, like a hard disk, which can temporarily slow down the program's execution.

Advantages of Demand Paging:

1. **Optimized Memory Usage:** Demand Paging minimizes the initial loading time by loading only the necessary parts of a program into memory when needed, thereby optimizing memory usage.
2. **Support for Larger Programs:** It allows the execution of larger programs or multiple programs concurrently, as it doesn't require the entire program to be loaded into memory at once.
3. **Reduced Initial Loading Time:** By loading only the essential pages initially, the system starts faster, as it doesn't need to load the entire program into memory before execution begins.
4. **Efficient Use of Resources:** Demand Paging efficiently utilizes available memory by bringing in only required parts of the program, allowing other programs to use memory space more effectively.

Disadvantages of Demand Paging:

1. **Potential for Page Faults:** When a required page is not in memory, it results in a page fault, leading to additional time spent fetching the page from secondary storage (like a hard disk), causing a temporary slowdown in program execution.
2. **Increased Overhead:** Managing and handling page faults, including the retrieval of pages from secondary storage, introduces additional overhead and may impact system performance.
3. **Complexity in Implementation:** Implementing demand paging requires sophisticated algorithms for page replacement, eviction policies, and efficient handling of page faults, which adds complexity to the system.
4. **Performance Impact:** Excessive page faults, particularly when the system's memory is heavily loaded or when there's a high demand for memory, can lead to increased disk access and degrade overall system performance.

How Paging is different from Demand Paging.

1. Basic Concept:

- **Paging:** It divides physical memory and processes into fixed-size blocks (pages).
- **Demand Paging:** It's a method of loading pages into memory only when they're needed.

2. Memory Usage:

- **Paging:** All pages of a process are loaded into memory before execution begins.
- **Demand Paging:** Only necessary pages are loaded initially; others are fetched on demand.

3. Loading Strategy:

- **Paging:** Loads entire program into memory before execution.

- **Demand Paging:** Loads only essential pages initially, fetching others as needed during execution.

4. Initial Loading Time:

- **Paging:** Longer initial loading time as entire program loads into memory.
- **Demand Paging:** Shorter initial loading time as only required pages load first.

5. Efficiency:

- **Paging:** May waste memory if the entire program isn't fully used.
- **Demand Paging:** Utilizes memory more efficiently by loading pages as needed.

6. System Start-Up:

- **Paging:** Slower system start as it loads full programs into memory.
- **Demand Paging:** Faster system start since only critical pages load initially.

7. Page Faults:

- **Paging:** Less susceptible to page faults as all pages are pre-loaded.
- **Demand Paging:** More likely to have page faults if required pages aren't in memory.

8. Resource Utilization:

- **Paging:** May not efficiently use memory if full program isn't utilized.
- **Demand Paging:** Efficiently uses memory by loading only necessary pages.

What is segmentation. Explain How address translation is done with a suitable example.

Segmentation is a memory management technique that divides a process's logical address space into variable-sized segments based on the program's structure or functionality (such as code, data, stack). Each segment represents a different part of the program and may vary in size, allowing for a more flexible memory allocation compared to fixed-size partitions in paging.

Address translation is the process of converting a logical address into its corresponding physical address using a segment table. Address translation in segmentation is done using a combination of a segment table and an offset within the segment. Here's an example:

Suppose a program has three segments: Code (segment 0), Data (segment 1), and Stack (segment 2). The segment table holds base addresses for each segment. For simplicity, let's assume:

- Segment 0 (Code) starts at address 4000.
- Segment 1 (Data) starts at address 6000.
- Segment 2 (Stack) starts at address 8000.

And let's say the program generates an instruction like: "Access address 350 in the Data segment."

1. Segmentation Process:

- The logical address consists of two parts: Segment Number and Offset.
- For the address 350 in the Data segment:
 - Segment Number: 1 (Data segment)
 - Offset: 350

2. Address Translation:

- The segment number (1) is used as an index to locate the base address of the Data segment in the segment table.
- The base address for the Data segment is 6000.

- Adding the offset (350) to the base address (6000) gives the physical memory address:
 - $\text{Physical Address} = \text{Base Address} + \text{Offset} = 6000 + 350 = 6350.$

Therefore, the logical address 350 in the Data segment translates to the physical address 6350 in memory, allowing the processor to access the required data.

Segmentation enables a more flexible memory management system by breaking down a program into different functional units (segments) and providing varying segment sizes.

What is the role of memory management in computer system ?

Memory management in a computer system is crucial as it performs several essential roles that ensure efficient utilization of the system's memory resources:

1. **Resource Allocation:** It allocates and assigns memory space to different processes or programs running concurrently, ensuring each gets the required memory for execution.
2. **Address Translation:** Converts logical addresses used by programs into physical addresses in memory, facilitating data access and execution by the CPU.
3. **Memory Protection:** Controls and safeguards memory areas, preventing unauthorized access and ensuring that processes cannot interfere with each other's memory space.
4. **Memory Sharing:** Facilitates sharing of memory segments or data among processes when required, enabling efficient communication and collaboration.
5. **Memory Optimization:** Manages memory efficiently to minimize fragmentation, ensuring the available memory is used optimally and preventing wastage.
6. **Virtual Memory Management:** Handles virtual memory systems, allowing the system to run more programs than the physical memory can accommodate by using secondary storage as an extension.

Write short notes on simplest memory management scheme.

The simplest memory management scheme is "**Single Partition Allocation**". (Now explain it...)

Explain Bare Machine ?

A "bare machine" refers to a computer system without any operating system or software installed. It represents the raw hardware components, including the processor, memory, storage devices, input/output interfaces, and other essential hardware elements.

In this state, when a computer is powered on, it lacks the necessary software layer to manage resources, handle user commands, or execute programs. Without an operating system, the machine is unable to perform complex tasks or run software applications.

In such a scenario, users interact with the hardware directly, issuing low-level commands or instructions to control the hardware components. Without the abstraction and management provided by an operating system, users would need to manage hardware resources manually, which is complex and highly inefficient.

A bare machine requires an operating system to enable users to interact with the hardware, manage resources, execute programs, and perform various tasks effectively. The operating system acts as an intermediary between the hardware and the user, providing a user-friendly interface and managing the system's resources to facilitate the execution of software applications.

Explain Resident Monitor ?

A Resident Monitor, in the context of early computing systems, refers to a part of the operating system that remains in memory (resident) throughout the entire duration of computer operation. It is a key component of the operating system that remains loaded in memory once the system is booted up and persists there until the system is powered off.

The primary function of the Resident Monitor is to oversee and manage various critical aspects of the computer system, such as handling input/output operations, managing memory, controlling the execution of programs, and providing essential system services. It acts as an interface between the hardware and user programs, ensuring proper execution and resource management.

Having the Resident Monitor reside in memory enables quick access to crucial operating system functionalities without the need to load it from external storage each time a system call or user program requires its services. This allows for more efficient and faster execution of programs as the essential system services are readily available in memory.