# 1. What is paging?  Explain with a paging hardware system with TLB.

Paging is a memory management scheme used in operating systems to allow processes to have a <u>non-contiguous memory allocation,</u> solving the issue of <u>memory fragmentation.</u> It divides both logical and physical memory into fixed-sized blocks, where:
- **Logical memory blocks are called pages.**
- **Physical memory blocks are called frames.**

Each process's pages are mapped to available frames in physical memory using a page table, which ensures <u>efficient memory utilization.</u>

## Key Features of Paging

1. <u>Eliminates External Fragmentation:</u>
   - Memory is divided into equal-sized frames, so it can be filled without gaps.

2. <u>Efficient Utilization of Memory:</u>
   - Processes are loaded into available frames, optimizing memory usage.

3. <u>Supports Virtual Memory:</u>
   - Allows processes larger than the available physical memory to execute.

## Paging Hardware System

Paging requires hardware support to translate logical addresses into physical addresses. The hardware components involved are:

1. **Page Table:**
   - Maintains a mapping of page numbers to frame numbers.
   - Each process has its own page table.

2. **Translation Lookaside Buffer (TLB):**
   - A high-speed cache storing a subset of the page table for faster access.
   - Helps reduce the time required for address translation.

## Working of Paging with TLB

1. **Logical Address Division:**
   - The CPU generates a logical address consisting of:
   - **Page number (p):** Index in the page table.
   - **Page offset (d):** Displacement within the page.

**2.     TLB Lookup:**
        •       The page number (p) is searched in the TLB.
        •       If found (TLB hit), the corresponding frame number is fetched, and the physical address is generated directly.
**3.     Page Table Lookup (if TLB Miss):**
        •       If not found in the TLB (TLB miss), the page table in memory is accessed to retrieve the frame number.

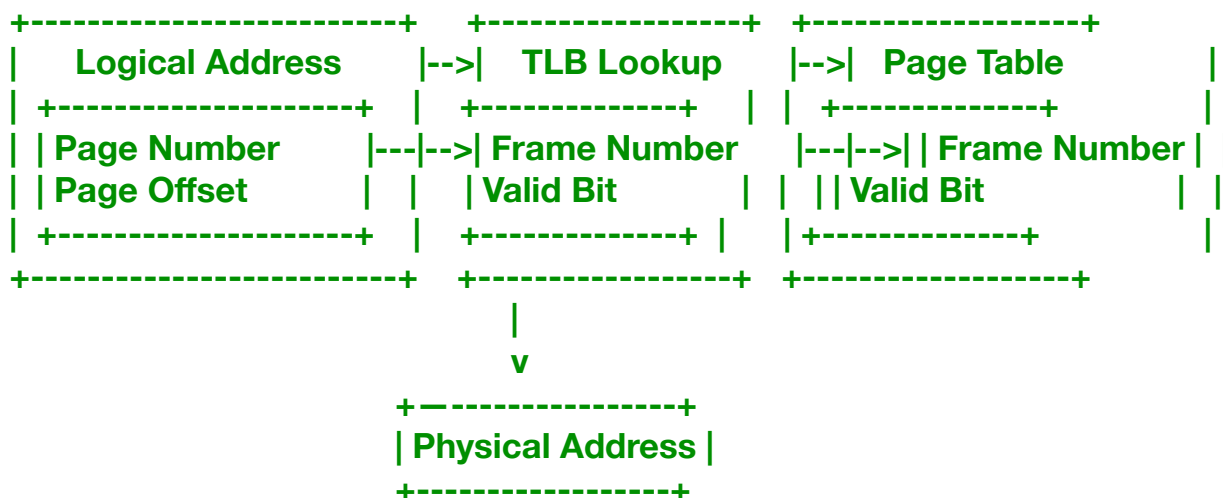**4.     Address Translation:**
        •       The frame number and page offset (d) are combined to generate the physical address.

**5.     TLB Update:**
        •       If the page table is accessed due to a TLB miss, the new mapping is added to the TLB for future reference.

## Advantages of Using TLB

        •       Faster Address Translation:
        •       Reduces the overhead of repeatedly accessing the page table.
        •       Improved System Performance:
        •       Enhances execution speed by reducing memory access times.

```
+--------------------------+   +---------------+   +-----------------+
|   Logical Address    |-->|   TLB Lookup   |-->|   Page Table       |
| +--------------------+  |   +-------------+   | | +-------------+         |
| | Page Number      |---|-->| Frame Number   |---|-->|| Frame Number |  |
| | Page Offset      |  |   | Valid Bit        |  |   ||| Valid Bit       |  |
| +--------------------+  |   +-------------+  |   | +-------------+        |
+--------------------------+   +---------------+   +-----------------+
                                      |
                                      v
                        +—---------------+
                        | Physical Address |
                        +-----------------+
```

# 2. With the help of a neat diagram, explain the various steps of address binding.

Address binding refers to the process of mapping <u>symbolic memory addresses</u> in a program to <u>physical memory addresses</u> in the computer system. This is essential because user programs typically refer to memory locations using symbolic names, which must be resolved to actual physical addresses before execution.

## Key Stages of Address Binding

There are **three stages** during which address binding can occur, depending on when the memory location is determined:

1. **Compile Time:**
   • If the memory location is known at compile time, absolute addresses are generated.
   • *Example:* The starting address is fixed, such as 2000.
   • *Limitation:* The program must always execute at the same memory location.

2. **Load Time:**
   • If memory location is unknown during compilation, relocatable addresses are generated.
   • When the program is loaded into memory, the loader translates these addresses to absolute addresses.

3. **Execution Time:**
   • If a program can move around in memory during execution, binding is delayed until runtime.
   • Requires hardware support (e.g., Memory Management Unit - MMU).
   • *Example:* Logical addresses are dynamically mapped to physical addresses.

## Logical vs Physical Addresses

   • **Logical Address:** The address generated by the CPU (used by a program).

   • **Physical Address:** The actual location in physical memory where the data resides.

# Steps of Address Binding with Diagram

**1.    Symbolic Address:**
- Initially, memory locations are referenced using <u>symbolic names.</u>
- **Example:** int x is stored at symbolic_address_x.

**2.    Compilation:**
- The compiler converts symbolic addresses into <u>relocatable addresses.</u>
- **Example:** x is stored at offset 20 within a segment.

**3.    Loading:**
- During loading, the <u>loader</u> converts <u>relocatable addresses into absolute addresses</u>, which are determined by the starting location of the segment in physical memory.

**4.    Execution:**
- At runtime, the CPU generates logical addresses, which are translated to physical addresses using the base register and limit register.

## Importance of Address Binding

1.    Provides abstraction, enabling user programs to work independently of the underlying hardware.

2.    Enables relocation, ensuring efficient use of memory.

3.    Facilitates dynamic allocation and sharing of resources in a multitasking environment.

```
+-------------------+        +----------------+        +-----------------+
| Compile-time Bind|------->|  Load-time Bind |———>| Execution-time Bind|
| (Source Code)     |        | (Loading in mem)|        | (Dynamic mapping) |
+-------------------+        +----------------+        +-----------------+
```

# 3. What is demand paging? Explain the steps in handling page faults using the appropriate diagram.

Demand paging is a <u>memory management technique</u> where pages of a process are loaded into memory only when they are required during execution. This minimizes the <u>memory footprint</u> of a process, allowing multiple programs to execute efficiently in limited physical memory.

## Key Features of Demand Paging

1. **Lazy Loading:**
   - Pages are not loaded into memory at the start of a process.
   - A page is brought into memory only when accessed.

2. **Reduction in Memory Usage:**
   - Only the required pages are loaded, saving memory resources.

3. **Supports Virtual Memory:**
   - Processes can have logical address spaces larger than physical memory.

4. **Improves Multiprogramming:**
   - Enables multiple processes to share limited memory by loading only required pages.

## Steps in Handling Page Faults

When a program references a page not currently in memory, a page fault occurs. The operating system must handle this fault to continue the execution.

## Steps:

1. **Memory Access Check:**
   - The memory address requested by the process is checked to ensure it is valid.

2. **Page Validity:**
   - If the address references an invalid page, the process is terminated.
   - If the page is valid but not in memory, a page fault is triggered.

3. **Locate Free Frame:**
   - The OS locates a free frame in physical memory to store the page.
   - If no free frame is available, a page replacement algorithm is used to evict a page from memory.

## 4. Load Page:

• The required page is fetched from secondary storage (e.g., disk) into the allocated frame.

## 5. Update Page Table:

• The page table entry is updated to mark the page as valid and record its frame number.

## 6. Restart Instruction:

• The instruction causing the page fault is restarted after the page is successfully loaded.

### Advantages of Demand Paging

## 1. Efficient Memory Usage:
• Reduces memory consumption by loading pages on demand.

## 2. Faster Process Startup:
• Only a minimal portion of the program is loaded initially.

## 3. Supports Large Applications:
• Applications larger than physical memory can execute seamlessly.

## Example of Effective Access Time (EAT):
• Memory access time: 200 nanoseconds.
• Average page fault service time: 8 milliseconds.
• Effective access time formula:

**EAT = (1 - P) \times MA + P \times (Page\ Fault\ Service\ Time)**

Where P is the probability of a page fault.

## Conclusion:

Demand paging optimizes memory usage and allows efficient multitasking in modern operating systems. By handling page faults efficiently, the system ensures smooth execution of processes even with limited physical memory.

## Diagram

```
+----------------------+      +----------------+      +------------------+      +-----------------+
|  CPU Generates Address| --->|   Check Page  |--->| Page Fault Occurs |—>  |Fetch Page from|
|  (Logical Address)    |     |  Table (Valid?) |    |                   |     |      Disk       |
+----------------------+      +----------------+      +------------------+      +-----------------+
                                      |                                                 |
                                      v                                                 v
                            +------------------+                            +-------------------+
                            | Update Page Table | <————————|Update TLB Cache  |
                            +------------------+                            +------------------+
                                      |
                                      v
                            +-------------------+
                            | Physical Address  |
                            +-------------------+
```

# 4. What is Segmentation? Explain the Basic Method of Segmentation with an Example

Segmentation is a memory management scheme that divides a program into logically distinct segments such as code, data, stack, etc., based on the user's view of the program. Each segment has a name and a length.

## Key Features of Segmentation

1. **Logical Addressing**:
   - Logical addresses consist of two parts:
   - Segment number (s): Identifies the segment.
   - Offset (d): Specifies the location within the segment.

2. **Segment Table:**
   - Maps each segment to its base address in physical memory.
   - Each table entry contains:
   - Base: Starting physical address of the segment.
   - Limit: Length of the segment.

3. **Protection and Isolation:**
   - Segments provide isolation and can be protected individually.

## Basic Method of Segmentation

1. **Compilation:**
   - Programs are divided into segments by the compiler.

   - **Common segments include:**

   - Code: Instructions for execution.
   - Global Variables: Data shared by all processes.
   - Heap: Memory dynamically allocated during runtime.
   - Stack: Memory used for function calls and local variables.

2. **Address Mapping:**
   - The segment table is used for mapping logical to physical addresses.
   - If the offset (d) exceeds the segment limit, a trap is generated to handle the error.

## 3.  Execution:
- During execution, the segment number is used to index the segment table.
- The physical address is calculated as:

$$\text{Physical Address} = \text{Base Address} + \text{Offset}$$

## Example:

Suppose a process has the following segments:
1. *Code:* 500 bytes starting at address 1000.
2. *Data:* 400 bytes starting at address 2000.
3. *Stack:* 300 bytes starting at address 3000.

Logical address (Segment 1, Offset 50):
- Base: 2000 (from segment table).
- Physical Address: 2000 + 50 = 2050.

```
+----------------------+        +------------------+              +------------------+
| Logical Address      |--->    | Segment Table    |     ----->| Physical Address   |
| +------------------+ |        |   +------------+  |              | +---------------+   |
| | Segment Number   | |        | | Base Address | |              | | Segment Data   |  |
| | Offset       | |  |        | | Limit Address| |              | | Segment Location|  |
| +------------------+ |        |   +------------+  |              |+---------------+    |
+----------------------+        +------------------+              +————————————+.
```

# 5. Discuss the structure of the page table with a suitable diagram.

A page table is a data structure used by the operating system to store the <u>mapping between logical page numbers and physical frame numbers in a paging system.</u>

## Structure:

1. **Single-Level Page Table:**
   - A single-level table maps each page to a frame.
   - **Entries include:**
     - <u>Page number:</u> Logical page index.
     - <u>Frame number:</u> Corresponding physical frame.
     - <u>Protection bits:</u> Access permissions for the page.

2. **Multi-Level Page Table:**

   - Used for large address spaces to reduce the size of a single page table.
   - Divides the logical address into multiple fields:
   - Outer page table points to inner page tables.
   - Inner page table maps pages to frames.

3. **Inverted Page Table:**

   - Contains one entry for each physical frame, mapping it to the owning process and page.

```
+-----------------+        +----------------+        +----------------+
| Page Number 1  |.  -->| Frame Number 1|-->| Valid/Invalid Bit|
+-----------------+        +----------------+        +----------------+
| Page Number 2  |.  -->| Frame Number 2]|—>| Valid/Invalid Bit|
+-----------------+        +----------------+        +-----------------+
| Page Number 3  |.  -->| Frame Number 3|-->| Valid/Invalid Bit|
+-----------------+        +----------------+        +-----------------+
```

# 6. What are Commonly used strategies   to select free hole from common holes

When allocating memory, the operating system manages a list of free holes and uses allocation strategies to choose a suitable hole for a process.

## Strategies:

1. **First Fit:**
   - Allocates the first hole large enough for the process.
   - **Advantages:** Simple and fast.
   - **Disadvantages:** May lead to fragmentation.

2. **Best Fit:**
   - Allocates the smallest hole that fits the process.
   - Searches the entire list unless ordered by size.
   - **Advantages:** Minimizes leftover space.
   - **Disadvantages:** May leave small unusable holes (external fragmentation).

3. **Worst Fit:**
   - Allocates the largest hole available.
   - **Advantages:** Leaves larger leftover holes.
   - **Disadvantages:** Often results in poor memory utilization.

```
+---------------------+
|Free Memory Holes  |
| +-----------------+ |
|| Hole 1 (Size 100) ||
| +-----------------+ |
|| Hole 2 (Size 200) ||
| +-----------------+ |
|| Hole 3 (Size 300) ||
| +-----------------+ |
+---------------------+
     |
     v
+-------------------------+   +-------------------------+   +--------------------------+
| First Fit (First hole fit)  |   | Best Fit (Smallest hole).  |   | Worst Fit (Largest hole)  |
+-------------------------+   +-------------------------+   +--------------------------+
```

# 7. What is address binding ? With neat diagrams explain steps involved in it ?

Address binding is the process of <u>mapping logical addresses</u> (used by programs) to <u>physical addresses</u> (used by the computer hardware). It is crucial for the execution of programs in memory and occurs in several stages: compile-time, load-time, and runtime. The binding ensures that the code and data are assigned correct memory locations during the execution of a program.

## Steps Involved in Address Binding:

## 1.    Compile-Time Address Binding:

•       **Description:** During the compilation process, the compiler generates the logical addresses. The program is compiled assuming that it will be loaded at a specific memory location.
•       **Process:** If a program's starting address is fixed (i.e., the compiler knows the memory location), address binding occurs at compile time.
•       **Outcome:** The logical addresses in the source code are directly mapped to physical addresses.

**Example:**
•       If a program is compiled to start at memory address 1000, it will use addresses 1000, 1001, 1002, etc., regardless of where the program is loaded in the memory.

## 2.    Load-Time Address Binding:

•       **Description:** During load-time, the program is loaded into memory, and the loader performs address binding.
•       **Process:** The program's logical addresses are mapped to physical addresses at the time the program is loaded into memory, allowing for relocation (e.g., the program could be loaded at address 2000 instead of 1000).
•       **Outcome:** The addresses generated by the program are modified according to the actual location in memory.
**Example:**
•       If a program is loaded at memory address 2000, all addresses in the program are incremented by 1000 to adjust for the difference between the logical and physical addresses.

## 3.    Runtime Address Binding:

•       **Description:** In some systems, the binding occurs during runtime. This is common in systems where the program can be dynamically relocated in memory.

• **Process:** The address translation is performed by the hardware (using a memory management unit, or MMU) during the program's execution, allowing the program to reference logical addresses that are translated to physical addresses on the fly.

• **Outcome:** The logical address of the program is dynamically mapped to the physical address during execution.
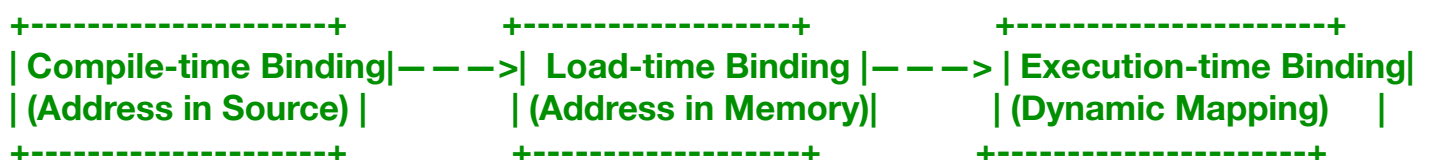
**Example:**

• In systems using virtual memory, the logical address generated by the CPU is translated to a physical address by the MMU at runtime.

## Summary of Address Binding:

• **Compile-time:** Binding occurs when the program is compiled, assuming a fixed memory location.

• **Load-time:** Binding occurs when the program is loaded into memory, with the memory address adjusted at load time.

• **Runtime:** Binding occurs during the execution of the program, allowing for dynamic relocation.

Each of these types of address binding enables flexibility in how programs are placed into memory, allowing for efficient memory utilization and process execution.

```
+-------------------+          +-----------------+          +--------------------+
| Compile-time Binding|———>|  Load-time Binding |———> | Execution-time Binding|
| (Address in Source) |          | (Address in Memory)|          | (Dynamic Mapping)    |
+-------------------+          +-----------------+          +--------------------+
```

# 8. Illustrate how demand paging effects system faults

**Demand paging**

Demand paging is a memory management technique in which pages of a process are only loaded into physical memory when they are needed, rather than being loaded all at once. This is in contrast to pure paging, where all pages of a process are loaded into memory during the process's execution.

**Page Fault:**

A page fault occurs when a process tries to access a page that is not currently loaded into physical memory. When this happens, the operating system must retrieve the page from secondary storage (like disk) and load it into memory.

## How Demand Paging Affects System Faults:

1. **Page Faults and Disk I/O:**
   • When a page is not in memory (a page fault), the system must load the page from the disk.
   • This leads to increased disk I/O, which can slow down the system if the fault rate is high.

2. **Increasing Latency:**
   • Latency increases because the system must access secondary storage (disk) every time a page fault occurs. This increases the time it takes to retrieve the required page and causes delays in the execution of programs.
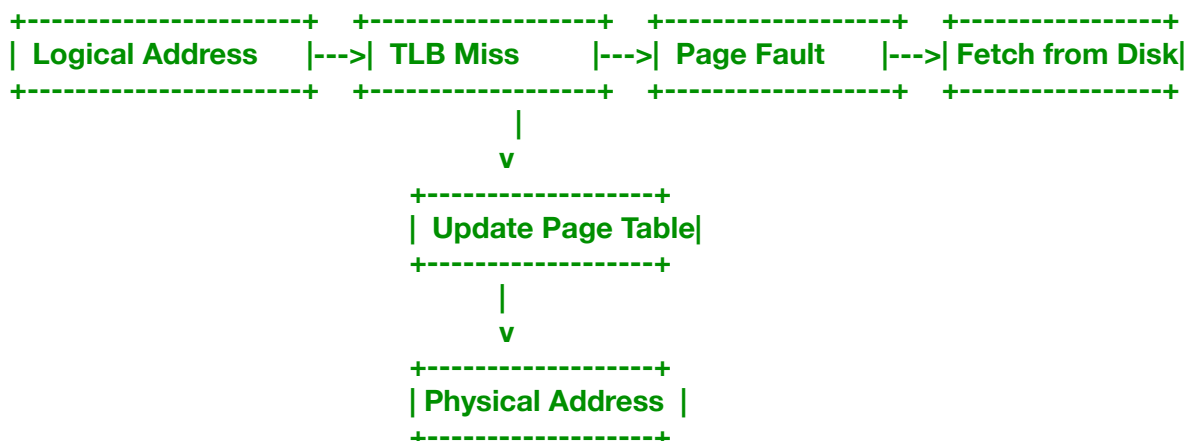
3. **Thrashing (Related Concept):**
   • If the system is constantly paging in and out of memory due to frequent page faults, this can lead to thrashing, where the system spends more time managing memory than executing the process.

Illustration of Demand Paging and Page Faults:

**Summary:**
   • Demand Paging reduces memory usage but can lead to frequent page faults, increasing disk I/O and latency.

```
+---------------------+   +-----------------+   +-----------------+   +----------------+
| Logical Address     |--->| TLB Miss        |--->| Page Fault      |--->| Fetch from Disk|
+---------------------+   +-----------------+   +-----------------+   +----------------+
                                  |
                                  v
                         +------------------+
                         | Update Page Table|
                         +------------------+
                                  |
                                  v
                         +------------------+
                         | Physical Address |
                         +------------------+
```

# What's trashing ? How it can be controlled

Trashing occurs when the system <u>spends most of its time swapping pages</u> <u>in and out</u> of memory due to a <u>high rate of page faults</u>. This occurs when there is insufficient physical memory for the active working set of processes, causing excessive paging and degrading overall system performance.

## Example of Trashing:
•      If a system is running multiple processes, each requiring more pages than can fit into memory, the system starts paging out the current pages and loading new pages constantly. This constant swapping leads to thrashing, where the CPU is kept busy managing memory, and actual execution slows down drastically.

## How Trashing Can Be Controlled:

1.    **Increase Physical Memory:**
•      One of the simplest ways to reduce thrashing is to add more RAM. More memory means the system can hold more pages, reducing the frequency of page faults and swaps.

2.    **Effective Page Replacement Algorithms:**
•      Implementing efficient page replacement algorithms (e.g., LRU – Least Recently Used) ensures that pages that are least likely to be used again are swapped out, reducing the chance of thrashing.

3.    **Process Scheduling:**
•      Limit the number of processes running concurrently. If too many processes are competing for memory, the system is more likely to thrash. By controlling the number of active processes, the system can better allocate memory to each process.

4.    **Working Set Model:**
•      The working set model ensures that each process has enough pages in memory to minimize page faults. This technique dynamically adjusts the memory allocation for processes based on their current needs.

5.    **Use of Virtual Memory:**
•      Efficient management of virtual memory can help control thrashing. The system must ensure that virtual memory pages are only swapped in when necessary and that processes are not overburdened with unnecessary paging.

## Illustration of Demand Paging and Page Faults:

---

1. Process attempts to access a page.
2. If the page is in memory:
   - No fault occurs.
   - The process continues as usual.

3. If the page is **not** in memory:
   - A **page fault** occurs.
   - The operating system must load the page from disk to memory.
   - The system may replace a different page in memory if necessary.

4. The program continues execution after the page has been loaded into memory.

---

**Conclusion:-**
   • Trashing is the result of excessive paging, causing the system to spend more time swapping pages than executing programs.
   • Trashing can be controlled by increasing physical memory, using efficient page replacement algorithms, limiting the number of processes, and managing virtual memory more effectively.

```
+----------------------------+        +— — — — — — — —+
| CPU generates requests  — —>| Frequent Page   |
| (Logical Address)          |        | Faults Occur    |
+----------------------------+        +------------------+
                                               |
                                               v
                                   +----------------------+
                                   |  System Thrashing    |
                                   |  (Slow Performance)  |
                                   +----------------------+
                                               |
                                               v
                                   +----------------------+
                                   | Control Thrashing    |
                                   | (Increase RAM/Swap)|
                                   +----------------------+
```