# PAPER: Operating Systems        Course Code: IT-306

**Q#1: What is Convey Effect?**

**Ans:** The Convey Effect refers to the phenomenon where simply putting your thoughts into words can improve your understanding of them. This effect highlights the powerful connection between language and cognition, demonstrating how the act of articulation can shape and refine our thinking.

**Examples**:

- ✓ A student studying for an exam explains the key concepts to a friend, improving their own understanding in the process.
- ✓ A manager clarifies their project plan by explaining it to their team, leading to better collaboration and execution.

**Remember**: The Convey Effect is not just about explaining things to others; it's about using the power of language to enhance your own thinking. By actively articulating your thoughts, you can gain a deeper understanding of yourself and the world around you.

**Q#2: Explain different states of process?**

**Ans:** Processes, whether physical, digital, or conceptual, often exist in various states throughout their lifecycle. Understanding these states is crucial for managing and optimizing processes effectively. Here's a breakdown of the different states a process can be in:

1. **Initial State:**
   - The process is in its initial stage, ready to be triggered or activated.
   - All necessary resources and inputs are available.
2. **Running State:**
   - The process is actively executing its tasks or steps.
   - Resources are being consumed, and outputs are being generated.
3. **Suspended State:**
   - The process is temporarily paused or halted.
   - The process can be resumed from its current state later.
4. **Waiting State:**
   - The process is waiting for a specific event or condition to proceed.
   - The process remains inactive until the awaited event occurs.
5. **Completed State:**
   - The process has successfully finished its execution.
   - The process may be terminated or archived depending on its purpose.
6. **Failed State:**
   - This could be due to missing resources, invalid inputs, or unexpected events.
   - The process may require intervention or restart to address the failure.
7. **Terminated State:**
   - The process has been intentionally stopped or ended.
   - The process is no longer active and cannot be resumed.
8. **Aborted State:**
   - The process has been abruptly stopped or interrupted.
   - This could be due to system crashes, power failures, or external events.

*MUHAMMAD BILAL RAFIQ*

**Q#3: Difference between long-term scheduler and short-term scheduler**.

**Ans:** In the realm of operating systems, schedulers play a crucial role in managing the allocation of resources, particularly CPU time, among various processes. While both long-term and short-term schedulers contribute to this task, they differ significantly in their scope, objectives, and decision-making processes.

## Long-Term Scheduler (LTS):

**Focus**: Overall system performance and resource utilization.

**Scope**: Determines which processes are admitted into the system and the degree of multiprogramming.

**Decision Factors**: Process priority, memory requirements, resource usage history, and overall system load.

**Frequency**: Less frequent, typically invoked when a new process arrives or when the system load changes significantly.

**Objectives:**

- Maximize CPU utilization.
- Maintain a balanced mix of processes to optimize resource usage.
- Prevent overloading the system with too many processes.

## Short-Term Scheduler (STS):

**Focus**: Efficient CPU utilization and process execution.

**Scope**: Selects a process from the pool of ready processes to be allocated the CPU.

**Decision Factors**: Process priority, CPU burst time, I/O requirements, and the current state of the system.

**Frequency**: Very frequent, typically invoked every few milliseconds or whenever a process changes state (e.g., from running to waiting).

**Objectives**:

- Minimize CPU idle time.
- Ensure fairness in CPU allocation among ready processes.
- Meet deadlines for real-time processes.

MUHAMMAD BILAL RAFIQ

**Q#4: Explain deadlock condition?**

**Ans:** Deadlock, a dreaded state in concurrent systems, refers to a situation where a set of processes are blocked indefinitely, preventing any of them from making progress. It's like a traffic jam where cars are stuck, unable to move because they are all blocking each other's way.

To understand deadlock, let's consider the four necessary conditions that must hold simultaneously for it to occur:

1. **Mutual Exclusion**: Each resource is held by only one process at a time. No sharing is allowed.
2. **Hold and Wait**: A process holding at least one resource is waiting to acquire additional resources currently held by other processes.
3. **No Preemption**: Resources cannot be forcibly taken away from a process. They must be released by the process holding them.
4. **Circular Wait**: There exists a circular chain of processes, each waiting for a resource held by the next process in the chain.

If all four conditions are met, the processes involved are stuck in a deadlock, unable to proceed further.

**Q#5: Define dispatcher latency.**

**Ans:** Dispatcher latency, also known as context switch latency, refers to the time it takes for a CPU to switch from executing one process to another. It's the delay experienced between the moment a process becomes ready to run and the moment it actually starts running on the CPU.

This latency arises due to the various tasks involved in switching contexts, such as:

❖ **Saving the state of the currently running process**: This includes storing the process's registers, program counter, and other relevant information in memory.
❖ **Loading the state of the new process**: This involves retrieving the new process's state from memory and loading it into the CPU registers.
❖ **Updating internal data structures**: The operating system needs to update its internal data structures to reflect the change in the running process.

The duration of dispatcher latency depends on various factors, including:

❖ **Hardware architecture**: The complexity of the CPU and memory system can influence the time it takes to save and load process states.
❖ **Operating system design**: The efficiency of the context switching algorithm and the amount of information that needs to be saved and loaded can impact latency.
❖ **Process characteristics**: Processes with larger memory footprints or more complex state information will take longer to switch contexts.

MUHAMMAD BILAL RAFIQ

**Q#6: Define demand paging.**

**Ans**: Demand paging, a cornerstone of modern operating systems, is a memory management technique that loads pages of a process into main memory only when they are actually needed. This contrasts with contiguous allocation, where the entire process is loaded into memory at once, even if only a small portion of it is used.

Here's how demand paging works:

- ❖ **Process Address Space**: The process's virtual address space is divided into fixed-size blocks called pages.
- ❖ **Page Table**: A page table is maintained for each process, mapping virtual pages to physical memory frames.
- ❖ **Initial State**: Initially, most pages are marked as "not present" in the page table, meaning they are not loaded in memory.
- ❖ **Page Fault**: When the process tries to access a page that is not present, a page fault occurs.
- ❖ **Demand Paging:** The operating system then fetches the required page from secondary storage (e.g., disk) and loads it into a free memory frame.
- ❖ **Page Table Update**: The page table is updated to reflect the new mapping between the virtual page and the physical frame.
- ❖ **Process Resumes**: The process resumes execution, now with the required page loaded in memory.

**Q#7: What is the process address space?**

**Ans:** The process address space, also known as virtual address space, is a fundamental concept in operating systems that provides a logical view of memory for a process. It's a contiguous range of addresses that the process can use to access its data and instructions, independent of the actual physical memory layout.

Here's a breakdown of the key aspects of the process address space:

- ❖ **Virtual Memory**: The process address space is a key component of virtual memory, which allows processes to have a larger address space than the available physical memory.
- ❖ **Logical View**: The process address space provides a logical view of memory for the process, hiding the complexities of physical memory management from the programmer.
- ❖ **Segmentation and Paging**: The process address space can be organized using segmentation or paging techniques.
- ❖ **Protection**: The process address space provides a mechanism for memory protection, preventing processes from accessing each other's memory or the operating system's memory.
- ❖ **Sharing**: The process address space can be shared among multiple processes, allowing them to access common data structures or libraries. This
- ❖ **Address Translation**: The MMU plays a crucial role in translating virtual addresses to physical addresses.

MUHAMMAD BILAL RAFIQ

**Q#8: Difference between paging and segmentation.**

**Ans:** Paging and segmentation are two distinct memory management techniques used in operating systems to organize and allocate memory to processes. While both share the goal of efficient memory utilization, they differ significantly in their approach and characteristics.

**Paging**:

Divides the process address space into fixed-size blocks called pages.

Each page is mapped to a physical memory frame of the same size.

Simpler to implement and manage due to the uniform page size.

Less prone to external fragmentation but can suffer from internal fragmentation.

Offers limited protection as access control is applied at the page level.

Easier to share pages among processes.

**Segmentation**:

Divides the process address space into variable-size segments based on logical units (e.g., code, data, and heap).

Each segment can have different access permissions and protection levels.

More complex to implement and manage due to variable segment sizes.

Less prone to internal fragmentation but can suffer from external fragmentation.

Offers more granular protection as access control is applied at the segment level.

More complex to share segments among processes.


**Q#9: Define turn-around time**.

**Ans:** Turnaround time (TAT), a fundamental metric in operating systems, measures the total time it takes for a process to complete its execution, from the moment it is submitted to the system until it finishes. It provides insights into the overall responsiveness and efficiency of the system.

Here's how TAT is calculated:

```
Turnaround Time = Completion Time - Submission Time
```

Where:

- ❖ Completion Time: The time when the process finishes execution.
- ❖ Submission Time: The time when the process is submitted to the system.

*MUHAMMAD BILAL RAFIQ*

**Q#10: Define two operations of semaphore.**

**Ans:** Semaphores, a fundamental synchronization primitive in operating systems, provide a mechanism for controlling access to shared resources and ensuring data consistency in concurrent environments. They operate through two key operations:

**Wait (P) Operation**:

- ❖ Also known as "down" operation or "acquire" operation.
- ❖ Decrements the semaphore value by 1.
- ❖ If the semaphore value becomes negative, the process is blocked until the value becomes positive again.
- ❖ This operation ensures that only one process can access the shared resource at a time, preventing race conditions and data corruption.

**Signal (V) Operation:**

- ❖ Also known as "up" operation or "release" operation.
- ❖ Increments the semaphore value by 1.
- ❖ If there are processes waiting for the semaphore, one of them is unblocked and allowed to proceed.
- ❖ This operation allows other processes to access the shared resource when it becomes available.

**Q#11: What is difference between user mode and system mode?**

**Ans:** In the realm of operating systems, processes operate in two distinct modes: user mode and system mode. These modes define the level of access and privileges a process has to system resources and functionality. Understanding the differences between these modes is crucial for comprehending how operating systems manage processes and protect the system from unauthorized access.

**User Mode:**

- ❖ **Limited Privileges**: Processes running in user mode have restricted access to system resources and functionality. They cannot directly access hardware devices, modify critical system data, or perform privileged operations.
- ❖ **Protection Mechanism**: User mode acts as a protective layer, preventing malicious or buggy applications from harming the system or other processes.
- ❖ **Typical Operations**: User mode processes typically perform tasks related to user applications, such as running programs, displaying data, and interacting with input devices.
- ❖ **Example:** A web browser, a text editor, or a media player running on your computer operates in user mode.

**System Mode:**

- ❖ **Elevated Privileges**: Processes running in system mode have full access to system resources and functionality. They can access hardware devices, modify system data, and perform privileged operations.
- ❖ **Kernel Access**: System mode processes operate within the kernel, the core of the operating system, giving them direct control over system resources and functionality.
- ❖ **Critical Operations**: System mode processes handle critical tasks such as memory management, process scheduling, device drivers, and file system operations.
- ❖ **Example:** The operating system itself, device drivers, and system utilities run in system mode.

**Q#12: Explain dispatch latency and its causes and effects.**

**Ans**: Dispatch latency, also known as context switch latency, refers to the time it takes for a CPU to switch from executing one process to another. It's the delay experienced between the moment a process becomes ready to run and the moment it actually starts running on the CPU.

**Causes of Dispatch Latency:**

1. **Saving the State of the Current Process**:
   - ❖ This involves storing the process's registers, program counter, and other relevant information in memory.
   - ❖ The amount of information to save depends on the process's complexity and the operating system's design.
2. **Loading the State of the New Process**:
   - ❖ This involves retrieving the new process's state from memory and loading it into the CPU registers.
   - ❖ The time taken for this depends on the size of the process's state and the speed of the memory system.
3. **Updating Internal Data Structures**:
   - ❖ The operating system needs to update its internal data structures to reflect the change in the running process.
   - ❖ This includes updating the process table, scheduling queues, and other relevant data structures.

**Effects of Dispatch Latency:**

4. **Increased Response Time**:
   - ❖ Users may experience delays in getting responses from applications due to the time spent waiting for context switches.
   - ❖ This can be particularly noticeable in interactive applications where responsiveness is critical.

5. **Reduced Throughput**:
   - ❖ The overall number of tasks completed per unit time can decrease due to the overhead of context switching.
   - ❖ This can impact the overall performance of the system, especially in multi-tasking environments.

6. **Unfairness**:
   - ❖ Processes with shorter execution times may be unfairly disadvantaged compared to processes with longer execution times due to the overhead of frequent context switches.
   - ❖ This can lead to imbalances in resource allocation and affect the overall fairness of the system.

**Q#13: What is starvation? Explain with the help of two examples.**

**Ans**: Starvation, in the context of operating systems, refers to a situation where a process is indefinitely prevented from accessing resources it needs to make progress. It's like being stuck in a queue, waiting for your turn, but never actually getting served.

**Causes of Starvation**:

- ❖ **Priority Inversion**: A higher-priority process continuously acquires a resource needed by a lower-priority process, preventing the lower-priority process from ever getting its turn.
- ❖ **Resource Deprivation**: A process is consistently denied access to essential resources, such as CPU time or memory, due to other processes monopolizing those resources.
- ❖ **Deadlock**: A circular dependency among processes, where each process holds a resource needed by another, prevents any of the processes from making progress.

**Examples of Starvation**:

**Priority Inversion**:

- ❖ A high-priority printer process acquires a lock on the printer device.
- ❖ A low-priority text editor process needs to print a document but is blocked because the printer process holds the lock.
- ❖ The printer process continues to print large documents, never releasing the lock.
- ❖ The text editor process remains starved, indefinitely waiting for its turn to print.

**Resource Deprivation:**

- ❖ A memory-intensive application consumes a large portion of available memory.
- ❖ A low-memory process, such as a system utility, is unable to allocate enough memory to run due to the memory-intensive application's heavy usage.
- ❖ The system utility remains starved, unable to perform its essential tasks due to the lack of available memory.

**Q#14: What is PCB? What are different types of information stored in PCB about a process? At the time of context-switch what role PCB plays?**

**Ans:** The Process Control Block (PCB), also known as the Task Control Block (TCB), is a data structure that contains all the information needed by the operating system to manage a process. It acts as the central repository of a process's state, including its execution context, resource allocation, and scheduling information.

**Types of Information Stored in a PCB:**

1) **Process Identification:**
   - ❖ Process ID (PID): A unique identifier assigned to the process.
   - ❖ Process Name: A descriptive name for the process.
   - ❖ User ID (UID): The user who owns the process.
2) **Process State:**
   - ❖ Current state of the process (e.g., running, waiting, ready).
   - ❖ Program counter (PC): The memory address of the next instruction to be executed.
   - ❖ CPU registers: The values of the CPU registers at the time of context switch.
3) **Process Scheduling Information:**
   - ❖ Priority: The priority level of the process.
   - ❖ Scheduling queue: The queue in which the process is waiting for its turn to run.
   - ❖ Time slices: The amount of CPU time allocated to the process.
4) **Memory Management Information:**
   - ❖ Memory address space: The virtual memory allocated to the process.
   - ❖ Page table: A table that maps virtual memory addresses to physical memory addresses.
5) **Accounting Information:**
   - ❖ CPU time used by the process.
   - ❖ I/O statistics (e.g., number of bytes read/written).
   - ❖ Memory usage statistics.
6) **I/O Information:**
   - ❖ List of open files.
   - ❖ List of devices assigned to the process.

**Role of PCB in Context Switching:**

During a context switch, the operating system saves the state of the currently running process in its PCB and loads the state of the new process from its PCB. This allows the operating system to quickly switch between processes without losing track of their execution state.

MUHAMMAD BILAL RAFIQ

# PAPER: Operating Systems          Course Code: IT-306

**Q#15: Write one solution for classical Procedure-Consumer problem with only three shared variables (in, out and buffer)**

**Ans**: The classical producer-consumer problem involves two processes: a producer that generates data and a consumer that consumes it. The challenge is to ensure that the producer and consumer operate concurrently without data loss or corruption.

**Solution with Three Shared Variables**:

1. **Shared Variables**:
   - ❖ in: An integer variable indicating the index of the next empty slot in the buffer.
   - ❖ out: An integer variable indicating the index of the next full slot in the buffer.
   - ❖ buffer: An array of fixed size that stores the data items produced by the producer and consumed by the consumer.

2. **Producer Process**:
   - ❖ While there is space in the buffer (i.e., (in + 1) % buffer_size != out):
   - ❖ Produce a data item.
   - ❖ Store the data item in the buffer at index in.
   - ❖ Increment in (wrap around if necessary).
   - ❖ Signal the consumer that a new data item is available.

3. **Consumer Process:**
   - ❖ While there is data in the buffer (i.e., in != out):
   - ❖ Retrieve the data item from the buffer at index out.
   - ❖ Increment out (wrap around if necessary).
   - ❖ Consume the data item.
   - ❖ Signal the producer that a slot is available in the buffer.

4. **Synchronization**:
   - ❖ To ensure proper synchronization between the producer and consumer, we can use semaphores or murexes.
   - ❖ empty: A semaphore initialized with the number of empty slots in the buffer.
   - ❖ full: A semaphore initialized with 0.

5. **Producer Process**:
   - ❖ Wait on empty.
   - ❖ Produce a data item and store it in the buffer.
   - ❖ Signal full.

6. **Consumer Process**:
   - ❖ Wait on full.
   - ❖ Retrieve and consume the data item from the buffer.
   - ❖ Signal empty.

**Correctness**:

This solution ensures that the producer and consumer operate concurrently without data loss or corruption. The use of semaphores or murexes guarantees mutual exclusion, preventing race conditions and ensuring data consistency.

MUHAMMAD BILAL RAFIQ

# PAPER: Operating Systems                    Course Code: IT-306

**Q#16: Write code for wait() and signal() operations of counting semaphore.**

**Ans:**

```c
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
#include <semaphore.h>

// Structure to represent a counting semaphore
typedef struct {
    int value;
    sem_t mutex;
    sem_t waiting_processes;
} CountingSemaphore;

// Initialize a counting semaphore
CountingSemaphore* counting_semaphore_init(int initial_value) {
    CountingSemaphore* semaphore = malloc(sizeof(CountingSemaphore));
    semaphore->value = initial_value;
    sem_init(&semaphore->mutex, 0, 1);
    sem_init(&semaphore->waiting_processes, 0, 0);
    return semaphore;
}

// Wait operation
void counting_semaphore_wait(CountingSemaphore* semaphore) {
    sem_wait(&semaphore->mutex);
    if (semaphore->value > 0) {
        semaphore->value--;
    } else {
        sem_post(&semaphore->mutex);
        sem_wait(&semaphore->waiting_processes);
    }
}

// Signal operation
void counting_semaphore_signal(CountingSemaphore* semaphore) {
    sem_wait(&semaphore->mutex);
    semaphore->value++;
    if (semaphore->value == 1) {
        sem_post(&semaphore->waiting_processes);
    }
    sem_post(&semaphore->mutex);
}

// Destroy a counting semaphore
void counting_semaphore_destroy(CountingSemaphore* semaphore) {
    sem_destroy(&semaphore->mutex);
    sem_destroy(&semaphore->waiting_processes);
    free(semaphore);
}
```

**Explanation**:

❖ The code defines a CountingSemaphore structure that includes the semaphore value, a mutex for synchronization, and a semaphore for managing waiting processes.

❖ The counting_semaphore_init() function initializes a counting semaphore with a given initial value.

❖ The counting_semaphore_wait() function decrements the semaphore value if it's positive. Otherwise, it adds the current process to the waiting queue and waits for a signal.

❖ The counting_semaphore_signal() function increments the semaphore value. If there are waiting processes, it removes the first one from the queue and signals it, allowing it to proceed.

❖ The counting_semaphore_destroy() function destroys a counting semaphore and releases its associated resources.

MUHAMMAD BILAL RAFIQ

**Usage:**

```c
// Create a counting semaphore with an initial value of 5
CountingSemaphore* semaphore = counting_semaphore_init(5);

// Producer process
void producer() {
    while (1) {
        // Produce some data
        // ...

        // Acquire the semaphore
        counting_semaphore_wait(semaphore);

        // Add the data to a shared resource
        // ...

        // Release the semaphore
        counting_semaphore_signal(semaphore);
    }
}

// Consumer process
void consumer() {
    while (1) {
        // Acquire the semaphore
        counting_semaphore_wait(semaphore);

        // Remove data from the shared resource
        // ...

        // Release the semaphore
        counting_semaphore_signal(semaphore);
    }
}

// Start the producer and consumer processes
int main() {
    // Create producer and consumer threads
    pthread_t producer_thread, consumer_thread;
    pthread_create(&producer_thread, NULL, (void*)producer, NULL);
    pthread_create(&consumer_thread, NULL, (void*)consumer, NULL);

    // Wait for the threads to complete
    pthread_join(producer_thread, NULL);
    pthread_join(consumer_thread, NULL);

    // Destroy the semaphore
    counting_semaphore_destroy(semaphore);

    return 0;
}
```

**Note:**

This code uses the POSIX semaphore API for synchronization. Ensure that you have the appropriate libraries and headers installed for your operating system.

MUHAMMAD BILAL RAFIQ

**Q#17: What is difference between microkernel and layered operating system structures?**

**Ans**: Both microkernel and layered operating system structures are designed to manage system resources and provide services to applications. However, they differ significantly in their architecture and approach.

**Microkernel**:

- ❖ Size: Smaller footprint, containing only essential functionalities like memory management, process management, and inter-process communication (IPC).
- ❖ Modularity: Highly modular, with additional services like file systems and device drivers implemented as separate user-space processes.
- ❖ Security: Enhanced security due to smaller attack surface and isolation of services in user space.
- ❖ Flexibility: Easier to extend and customize by adding or removing modules.
- ❖ Performance: Potentially slower due to increased context switching between kernel and user space.
- ❖ Examples: MINIX, L4, QNX

**Layered Operating System**:

- ❖ Size: Larger footprint, with multiple layers stacked on top of each other, each providing specific functionalities.
- ❖ Modularity: Less modular, with services tightly integrated within the kernel.
- ❖ Security: Potentially less secure due to larger attack surface and increased complexity.
- ❖ Flexibility: Less flexible, as modifying the kernel requires more effort.
- ❖ Performance: Potentially faster due to reduced context switching within the kernel.
- ❖ Examples: Unix, Linux, Windows

**Q#18: In multithreaded application Many-to-One model is implemented. How does this application handle the threads?**

**Ans**: In a multithreaded application, the Many-to-One model implies that multiple threads are mapped to a single operating system (OS) thread. This approach offers several advantages and disadvantages compared to other models like One-to-One and Many-to-Many.

**Advantages**:

- ❖ Reduced overhead: Mapping multiple threads to a single OS thread minimizes context switching overhead, potentially improving performance.
- ❖ Simplified synchronization: Synchronization between threads becomes easier as they share the same execution context.

**Disadvantages**:

- ❖ Limited parallelism: All threads share the same OS thread, limiting true parallelism and potentially hindering performance on multi-core systems.
- ❖ Blocking issues: If one thread blocks, all other threads sharing the same OS thread are also blocked, potentially leading to performance bottlenecks.

**Q#19: Why every thread has its own stack?**

**Ans**: In a multithreaded application, each thread has its own stack for several crucial reasons:

1. **Isolation and Protection:**
   - ❖ Each thread's stack acts as a separate memory space, isolating its data from other threads. This prevents accidental data corruption or interference between threads, ensuring data integrity and program stability.
   - ❖ This isolation also protects threads from memory access errors or security vulnerabilities in other threads.

2. **Efficient Context Switching:**
   - ❖ When the operating system switches between threads, it needs to save the state of the current thread and load the state of the next thread. Having a separate stack for each thread allows for efficient context switching, as the OS only needs to save and restore the stack pointer.
   - ❖ This reduces the overhead associated with context switching, improving the overall performance of the multithreaded application.

3. **Local Variables and Function Calls:**
   - ❖ Each thread's stack stores its local variables and the return addresses for function calls. This allows each thread to maintain its own execution context and call functions independently without affecting other threads.
   - ❖ Without separate stacks, threads would share the same local variables and function call stack, leading to potential conflicts and unpredictable behavior.

4. **Exception Handling:**
   - ❖ When an exception occurs within a thread, the exception information is stored on the thread's stack. This allows the exception handler to access the necessary context and handle the exception appropriately.
   - ❖ If all threads shared the same stack, exception handling would become more complex and potentially unreliable.

5. **Thread-Local Storage:**
   - ❖ Each thread can have its own thread-local storage (TLS) associated with its stack. This allows threads to store private data that is accessible only to that specific thread.
   - ❖ TLS provides a convenient way for threads to store temporary data or maintain thread-specific settings without worrying about conflicts with other threads.

MUHAMMAD BILAL RAFIQ

# PAPER: Operating Systems                    Course Code: IT-306

**Q#20: What is zombie state of a process?**

**Ans**: The zombie state, also known as the "defunct" state, is a special state a process enters after it has terminated but its parent process has not yet performed a wait() or waitpid() system call to collect its exit status.

**How it happens:**

- ❖ **Process Termination**: A process terminates either normally by reaching its end or abnormally due to a signal or error.
- ❖ **Parent Process Responsibilities**: The parent process is responsible for collecting the exit status of its child process using wait() or waitpid().
- ❖ **Zombie State**: If the parent process does not perform this collection, the child process enters the zombie state.
- ❖ **Zombie Characteristics**: In the zombie state, the process:
  - No longer consumes CPU time or resources.
  - Still occupies an entry in the process table.
  - Has a process ID (PID) but cannot be interacted with.
  - Has an exit status that can be retrieved by the parent process.

**Why it exists:**

The zombie state exists for several reasons:

- ❖ Inform Parent Process: It allows the parent process to learn about the child's termination and retrieve its exit status.
- ❖ Resource Cleanup: The parent process can use the exit status to determine the reason for termination and perform any necessary cleanup tasks.
- ❖ Historical Information: The zombie state provides a record of the child's execution for debugging or monitoring purposes.

**Q#21: Name the inter-process communication tools used in UNIX operating system?**

**Ans:** The UNIX operating system offers a variety of tools and mechanisms for inter-process communication (IPC), allowing processes to exchange data and synchronize their actions. Here are some of the most commonly used IPC tools:

1. **Pipes**:
   - ❖ Type: Unidirectional, one-way communication between two related processes.
   - ❖ Mechanism: Creates a virtual pipe in memory, where one process writes data and the other reads it.
   - ❖ Use Cases: Simple data transfer between parent and child processes, command-line filtering (e.g., ls | grep).
2. **FIFOs (Named Pipes)**:
   - ❖ Type: Unidirectional or bidirectional, named pipes accessible by any process.
   - ❖ Mechanism: Similar to pipes, but data persists in the file system until read.

MUHAMMAD BILAL RAFIQ

❖ Use Cases: Communication between unrelated processes, data exchange across the network.

3. **Sockets:**
   ❖ Type: Bidirectional, network-based communication between processes on the same or different machines.
   ❖ Mechanism: Uses IP addresses and port numbers to establish connections.
   ❖ Use Cases: Network communication, client-server applications, distributed systems.

4. **Shared Memory:**
   ❖ Type: Bidirectional, high-speed data sharing between processes.
   ❖ Mechanism: Allocates a shared memory segment accessible by multiple processes.
   ❖ Use Cases: Efficient data exchange between processes with high data transfer rates.

5. **Semaphores:**
   ❖ Type: Synchronization mechanism for controlling access to shared resources.
   ❖ Mechanism: Uses a counter to regulate the number of processes accessing a resource.
   ❖ Use Cases: Mutual exclusion, resource allocation, process synchronization.

6. **Message Queues:**
   ❖ Type: Asynchronous, message-passing mechanism for communication between processes.
   ❖ Mechanism: Stores messages in a queue, allowing processes to send and receive messages independently.
   ❖ Use Cases: Loosely coupled communication, buffering data, asynchronous processing.

7. **Signals:**
   ❖ Type: Asynchronous notification mechanism to send signals between processes.
   ❖ Mechanism: Sends a signal to a process, interrupting its execution and potentially triggering a specific action.
   ❖ Use Cases: Event notification, process termination, inter-process control.

8. **Remote Procedure Calls (RPC):**
   ❖ Type: High-level mechanism for invoking procedures on remote machines.
   ❖ Mechanism: Transparently calls procedures on remote systems as if they were local.
   ❖ Use Cases: Distributed computing, client-server applications, network services.

**Q#22: What does a child process inherits from its parent process?**

**Ans**: When a child process is created by forking a parent process, it inherits certain attributes and resources from its parent. Here are some of the things that a child process typically inherits from its parent:

❖ **Process ID (PID):** The child process is assigned a unique PID, which is different from the parent's PID but is based on it.
❖ **Parent Process ID (PPID):** The child process is aware of its parent's PID and can access it using system calls. The PPID helps establish the parent-child relationship.
❖ **Open file descriptors:** The child process inherits the open file descriptors from the parent. It means that if the parent had files or network connections open, the child process will have access to them.

MUHAMMAD BILAL RAFIQ

❖ **Environment variables**: The child process inherits the environment variables of the parent process. These variables include information such as paths, system configuration, and other settings.

❖ **Working directory**: The child process starts with the same working directory as the parent process. It means that relative file paths will be resolved from the same location.

❖ **Signal handlers**: A child process inherits the signal handlers set by the parent. Signal handlers are used to handle events such as interrupts or termination signals.

❖ **User and group ID**: The child process inherits the user and group ID from its parent. This determines the permissions and access rights the child process has.

It's important to note that once the child process is created, it becomes an independent process and can modify or change the inherited attributes as needed.

**Q#23: What are the different techniques for the evaluation of scheduling algorithms?**

**Ans**: Evaluating the performance of scheduling algorithms is crucial to determine their effectiveness and suitability for different scenarios. Several techniques can be used to assess and compare scheduling algorithms, each with its own strengths and limitations.

1. **Simulation**:
   ❖ Concept: Simulating the execution of tasks on a virtual system with different scheduling algorithms.
   ❖ Advantages: Allows for controlled experimentation, testing various scenarios and workloads.
   ❖ Disadvantages: Can be time-consuming and may not accurately reflect real-world behavior.
2. **Analytical Modeling:**
   ❖ Concept: Using mathematical models to predict the performance of scheduling algorithms under specific assumptions.
   ❖ Advantages: Provides insights into theoretical behavior and can be computationally efficient.
   ❖ Disadvantages: May require simplifying assumptions and may not capture all aspects of real-world systems.
3. **Benchmarking:**
   ❖ Concept: Running standard benchmarks or workloads on different systems with different scheduling algorithms.
   ❖ Advantages: Provides a practical comparison of performance under realistic conditions.
   ❖ Disadvantages: May not be representative of all possible workloads and can be influenced by hardware and system configuration.
4. **Measurement and Tracing:**
   ❖ Concept: Collecting performance data from real-world systems running different scheduling algorithms.
   ❖ Advantages: Provides real-world performance insights and can identify bottlenecks and inefficiencies.

MUHAMMAD BILAL RAFIQ

❖ Disadvantages: Can be challenging to set up and may require specialized tools and expertise.

5. **User Perception:**
   ❖ Concept: Evaluating the perceived performance of scheduling algorithms based on user experience.
   ❖ Advantages: Captures the subjective experience of users and can be valuable for interactive systems.
   ❖ Disadvantages: Can be subjective and difficult to quantify, and may not reflect objective performance metrics.

6. **Hybrid Approaches:**
   ❖ Concept: Combining multiple techniques, such as simulation and measurement, to provide a more comprehensive evaluation.
   ❖ Advantages: Leverages the strengths of different techniques and can provide a more holistic understanding of performance.
   ❖ Disadvantages: Can be more complex and require additional resources and expertise.

**Q#24: Mention the three problems that may be caused by the wrong initialization or placement of wait ( ) and signal ( ) operations in the use of semaphores.**

**Ans:** Using semaphores for synchronization requires careful initialization and placement of wait() and signal() operations to avoid potential problems. Here are three common issues that can arise from incorrect usage:

1. **Deadlock:**
   ❖ Cause: A deadlock occurs when two or more processes are waiting for each other to release a resource they both need.
   ❖ Example: Process A needs resource X and Y, while process B needs resource Y and X. If process A acquires X and then waits for Y, while process B acquires Y and then waits for X, both processes will be stuck waiting indefinitely, creating a deadlock.
   ❖ Prevention: Ensure that processes acquire resources in a consistent order and release them in the reverse order. Avoid placing wait() operations before acquiring resources or signal() operations after releasing them.

2. **Starvation:**
   ❖ Cause: Starvation occurs when a process is repeatedly denied access to a resource due to other processes being prioritized.
   ❖ Example: A high-priority process repeatedly acquires a resource and signals it, preventing a low-priority process from ever acquiring the resource, leading to starvation.
   ❖ Prevention: Use fair scheduling algorithms to ensure all processes have a chance to acquire resources. Avoid holding resources for longer than necessary and use signal() operations judiciously to allow other processes access.

3. **Race Conditions:**
   - ❖ Cause: A race condition occurs when the order of execution of wait() and signal() operations is critical and can lead to unpredictable behavior.
   - ❖ Example: Two processes try to increment a shared counter. If one process increments the counter before the other process waits for it, the final value may be incorrect due to the race condition.
   - ❖ Prevention: Use semaphores in conjunction with other synchronization mechanisms like murexes or critical sections to ensure proper ordering of operations and avoid race conditions.

Q#25: What do you mean by a busy waiting semaphore or spin lock?

Ans: The terms "busy waiting semaphore" and "spin lock" are actually quite close, but not exactly the same. Here's a breakdown:

**Busy waiting:**

- ❖ This is a technique where a process or thread keeps checking a condition repeatedly in a loop.
- ❖ It's like waiting impatiently by constantly tapping your foot.
- ❖ While busy waiting, the process wastes CPU resources because it's not doing any useful work, just checking.

**Spin lock:**

- ❖ This is a type of synchronization mechanism that uses busy waiting.
- ❖ A spin lock is a flag that only one thread can hold at a time.
- ❖ Threads trying to acquire the lock keep checking (spinning) until it becomes available.
- ❖ Spin locks are good for short critical sections (shared resource access) where the wait time is expected to be low.

**Q#26: In a 64 bit machine, with 16GB RAM and 8KB page size. How many entries will there be in the page if it is an inverted page table?**

**Ans**:

**Given**:

64-bit machine

16GB RAM (16 * 2^30 bytes)

8KB page size (8 * 2^10 bytes)

**Calculation**:

1. **Number of Pages:** Divide the total RAM size by the page size:

```
Number of Pages = RAM Size / Page Size
Number of Pages = (16 * 2^30 bytes) / (8 * 2^10 bytes)
Number of Pages = 2^20
```

2. **Number of Entries in Inverted Page Table**: Since each page entry in the inverted page table corresponds to a frame in physical memory, the number of entries will be equal to the number of pages:

```
Number of Entries = Number of Pages
Number of Entries = 2^20
```

Therefore, in a 64-bit machine with 16GB RAM and 8KB page size, an inverted page table will have **2^20 entries**.

**Q#27: Explain the difference between layered and module based operating system structure?**

**Ans**: Both layered and module-based operating system structures aim to organize the operating system into smaller, manageable units. However, they differ in their approach and level of modularity.

**Layered Structure**:

❖ **Concept**: The operating system is divided into distinct layers, each with specific functionalities and well-defined interfaces.
❖ **Characteristics**:
  • Strict hierarchy: Each layer depends on the functionalities of the layer below it.
  • Limited modularity: Modifications to a layer often require changes in other layers due to their interdependence.
  • Examples: Unix, Linux, Windows

*MUHAMMAD BILAL RAFIQ*

**Module-Based Structure:**

- ❖ **Concept**: The operating system is composed of independent modules that communicate with each other through well-defined interfaces.
- ❖ **Characteristics**:
  - High modularity: Modules can be added, removed, or modified without affecting other modules.
  - Loose coupling: Modules interact through interfaces, minimizing dependencies and promoting flexibility.
  - Examples: MINIX, QNX

**Q#28: Describe the difference among short-term, medium term, and long term scheduling.**

**Ans:** Operating systems employ different scheduling algorithms at various levels to manage the execution of processes. These levels are categorized as short-term, medium-term, and long-term scheduling, each with distinct goals and responsibilities.

1. **Short-Term Scheduling (CPU Scheduling):**
   - ❖ **Goal**: Determines which process to allocate the CPU to at any given time.
   - ❖ **Timeframe**: Occurs on a millisecond or microsecond timescale.
   - ❖ **Algorithm Focus:** Prioritizes processes based on factors like CPU burst time, I/O requirements, and priority levels.
   - ❖ **Common Algorithms**: Round Robin, First-Come-First-Served, Shortest Job First, Priority Scheduling.
2. **Medium-Term Scheduling (Swapping):**
   - ❖ **Goal**: Decides which processes should be swapped out to secondary storage (e.g., disk) and which should be brought back into memory for execution.
   - ❖ **Timeframe**: Occurs on a timescale of seconds or minutes.
   - ❖ **Algorithm Focus**: Considers factors like memory availability, process priority, and resource utilization.
   - ❖ **Common Algorithms**: Least Recently Used (LRU), Most Recently Used (MRU), Second-Chance Algorithm.
3. **Long-Term Scheduling (Job Scheduling):**
   - ❖ **Goal**: Controls the degree of multiprogramming by determining which jobs should be admitted into the system and the level of resources allocated to them.
   - ❖ **Timeframe:** Occurs on a timescale of minutes or hours.
   - ❖ **Algorithm Focus**: Considers factors like job size, resource requirements, turnaround time, and priority.
   - ❖ **Common Algorithms**: First-Come-First-Served, Priority Scheduling, Round Robin.

MUHAMMAD BILAL RAFIQ

**Q#29: What are the two differences between kernel level thread and user level thread? Under what circumstance is one type better than other?**

**Ans**: Differences between Kernel-Level and User-Level Threads

1. **Management and Scheduling:**
   - ❖ **Kernel-Level Threads**: Managed and scheduled directly by the operating system kernel. The kernel has complete control over thread creation, scheduling, and synchronization.
   - ❖ **User-Level Threads**: Managed and scheduled by a user-level library or application. The operating system is unaware of these threads and only sees the process as a whole.
2. **Context Switching:**
   - ❖ **Kernel-Level Threads:** Context switching between kernel-level threads involves switching the entire kernel state, which can be expensive due to the need to save and restore kernel data structures.
   - ❖ **User-Level Threads**: Context switching between user-level threads is faster as it only involves switching the user-level context within the process, without affecting the kernel state.

**When to Choose Which Type:**

1. **Kernel-Level Threads**:
   - ❖ Advantages: More efficient for I/O-bound tasks due to faster context switching, better support for system calls and inter-process communication.
   - ❖ Disadvantages: More complex to implement and manage, requires privileged access, less portable across different operating systems.
   - ❖ Best for: Applications with frequent I/O operations, multi-core systems where efficient context switching is crucial.
2. **User-Level Threads**:
   - ❖ Advantages: Easier to implement and manage, more portable across different operating systems, allows for more flexible scheduling policies.
   - ❖ Disadvantages: Less efficient for I/O-bound tasks due to slower context switching, limited support for system calls and inter-process communication.
   - ❖ Best for: Applications with CPU-bound tasks, situations where portability and ease of implementation are prioritized.

MUHAMMAD BILAL RAFIQ

**Q#30: In a multithreaded application, if we implement Many-to-Many model then how does this application handle the threads?**

**Ans:** In a multithreaded application, the Many-to-Many model implies that multiple threads are mapped to multiple operating system (OS) threads. This approach offers several advantages and disadvantages compared to other models like One-to-One and Many-to-One.

**Advantages:**

- ❖ Increased parallelism: Multiple threads can execute concurrently on multiple CPU cores, potentially leading to significant performance improvements.
- ❖ Improved resource utilization: Threads can share resources like memory and file handles, potentially reducing resource consumption.

**Disadvantages:**

- ❖ Increased complexity: Managing multiple threads and their interactions becomes more complex, requiring careful synchronization and coordination.
- ❖ Overhead of context switching: Frequent context switching between threads can lead to performance overhead, especially on systems with limited resources.

**Q#31: Mention the three problems that may be caused by the wrong initialization or placement of wait ( ) and signal ( ) operations in the use of semaphores.**

**Ans**: Using semaphores for synchronization requires careful initialization and placement of wait() and signal() operations to avoid potential problems. Here are three common issues that can arise from incorrect usage:

1. **Deadlock:**
   - ❖ **Cause**: A deadlock occurs when two or more processes are waiting for each other to release a resource they both need.
   - ❖ **Example**: Process A needs resource X and Y, while process B needs resource Y and X. If process A acquires X and then waits for Y, while process B acquires Y and then waits for X, both processes will be stuck waiting indefinitely, creating a deadlock.
   - ❖ **Prevention**: Ensure that processes acquire resources in a consistent order and release them in the reverse order. Avoid placing wait() operations before acquiring resources or signal() operations after releasing them.
2. **Starvation:**
   - ❖ **Cause**: Starvation occurs when a process is repeatedly denied access to a resource due to other processes being prioritized.
   - ❖ **Example**: A high-priority process repeatedly acquires a resource and signals it, preventing a low-priority process from ever acquiring the resource, leading to starvation.
   - ❖ **Prevention**: Use fair scheduling algorithms to ensure all processes have a chance to acquire resources. Avoid holding resources for longer than necessary and use signal() operations judiciously to allow other processes access.

3. **Race Conditions:**
   - ❖ **Cause**: A race condition occurs when the order of execution of wait() and signal() operations is critical and can lead to unpredictable behavior.
   - ❖ **Example**: Two processes try to increment a shared counter. If one process increments the counter before the other process waits for it, the final value may be incorrect due to the race condition.
   - ❖ **Prevention**: Use semaphores in conjunction with other synchronization mechanisms like murexes or critical sections to ensure proper ordering of operations and avoid race conditions.

**Q#32: What are the main differences between operating systems for mainframe computers and personal computers?**

**Ans**: While both mainframe and personal computer (PC) operating systems share some fundamental functionalities, there are significant differences between them due to the distinct nature of their intended use cases and hardware environments. Here's a breakdown of the key distinctions:

1. **Purpose and Target Audience:**
   - ❖ Mainframe OS: Designed for high-volume transaction processing, large-scale data management, and mission-critical applications in enterprise environments.
   - ❖ PC OS: Primarily focused on individual user productivity, entertainment, and general-purpose computing tasks. Targets individual users and small businesses.
2. **Hardware Architecture:**
   - ❖ Mainframe OS: Runs on powerful, fault-tolerant mainframe computers with specialized hardware features for high availability and scalability.
   - ❖ PC OS: Designed for personal computers with diverse hardware configurations, ranging from low-end to high-performance desktops and laptops.
3. **User Interface:**
   - ❖ Mainframe OS: Traditionally command-line based, with graphical user interfaces (GUIs) becoming more prevalent in recent years.
   - ❖ PC OS: Primarily GUI-based, with intuitive interfaces designed for ease of use by individual users.
4. **Security:**
   - ❖ Mainframe OS: Emphasizes robust security features to protect sensitive data and ensure system integrity, often including multi-level access control and encryption.
   - ❖ PC OS: Security is important, but the focus is more on protecting individual user data and preventing malware infections.
5. **Scalability and Performance:**
   - ❖ Mainframe OS: Designed for high scalability and performance to handle large workloads and numerous concurrent users.
   - ❖ PC OS: Performance varies depending on hardware configuration, but generally less scalable than mainframe OS.

MUHAMMAD BILAL RAFIQ

6. **Cost:**
   - ❖ Mainframe OS: Typically more expensive due to the higher cost of mainframe hardware and specialized software licenses.
   - ❖ PC OS: Generally less expensive, with a wide range of options available at different price points.
7. **Applications:**
   - ❖ Mainframe OS: Supports enterprise-level applications for transaction processing, database management, data analytics, and mission-critical operations.
   - ❖ PC OS: Supports a wide variety of applications for individual productivity, entertainment, communication, and general-purpose computing.
8. **Open Source vs. Proprietary:**
   - ❖ Mainframe OS: Traditionally proprietary, but some open-source options are emerging.
   - ❖ PC OS: Both open-source (e.g., Linux) and proprietary (e.g., Windows, macOS) options are widely available.
9. **Support and Maintenance:**
   - ❖ Mainframe OS: Typically comes with dedicated support and maintenance services from the vendor.
   - ❖ PC OS: Support options vary depending on the specific OS and vendor, with both free and paid options available.
10. **Community and Resources:**
    - ❖ Mainframe OS: Smaller community compared to PC OS, but with specialized resources and expertise available.
    - ❖ PC OS: Large and active community with extensive resources, documentation, and support available online and through various channels.

**Q#33: What are the four necessary conditions for a deadlock to occur?**

**Ans:** A deadlock is a situation where two or more processes are blocked indefinitely, waiting for resources that are held by other processes. For a deadlock to occur, four necessary conditions must be met simultaneously:

1. **Mutual Exclusion**: Each resource can be held by only one process at a time. No other process can access the resource while it is being held.
2. **Hold and Wait**: A process holding at least one resource is waiting to acquire additional resources that are currently held by other processes.
3. **No Preemption**: Resources cannot be forcibly taken away from a process. They must be released by the process holding them.
4. **Circular Wait:** There exists a circular chain of processes, each waiting for a resource that is held by the next process in the chain.

**If any of these four conditions is not met**, a deadlock cannot occur. For example, if resources can be preempted, a process can be interrupted and its resources taken away, breaking the deadlock cycle. Similarly, if there is no circular wait, processes can eventually acquire the resources they need and avoid deadlock.

MUHAMMAD BILAL RAFIQ

**Q#34: Explain the difference between internal and external fragmentation.**

**Ans**: **Internal Fragmentation:**

- ❖ **Concept**: Occurs when allocated memory blocks are larger than the actual memory required by the process, resulting in unused space within the block.
- ❖ **Causes**:
    - Fixed-size allocation schemes (e.g., contiguous memory allocation).
    - Memory allocation algorithms that do not find the best fit for the process's memory requirements.
- ❖ **Example**: Allocating a 4KB memory block to a process that only needs 2KB, leaving 2KB of unused space within the block.
- ❖ **Impact**: Reduces memory utilization and can lead to memory exhaustion even when there is sufficient free memory available.

**External Fragmentation:**

- ❖ **Concept**: Occurs when free memory is broken into small, non-contiguous chunks that are too small to be used by any process. This fragmented free memory cannot be allocated to any process, even though the total amount of free memory may be sufficient.
- ❖ **Causes**:
    - Frequent allocation and deallocation of memory blocks.
    - Compaction algorithms that are not efficient in merging free memory blocks.
- ❖ **Example**: After several memory allocations and deallocations, the free memory is scattered into small chunks of 1KB, 2KB, and 5KB, which are too small to accommodate a process requiring 8KB.
- ❖ **Impact**: Reduces memory utilization and can lead to memory exhaustion even when there is enough total free memory.

**Q#35: What is the copy-on-write feature and under what circumstances is it beneficial to use this feature?**

**Ans**: Copy-on-write is a memory optimization technique that allows multiple processes to share the same pages of memory until one of them attempts to modify the data. At that point, a copy of the page is created for the modifying process, and the original page remains unchanged for the other processes.

**Benefits**:

- ❖ Reduced memory consumption: COW avoids unnecessary duplication of data, especially when multiple processes share read-only data. This can significantly reduce memory usage and improve system performance.
- ❖ Simplified memory management: COW eliminates the need for explicit copying of data when processes fork or share memory, simplifying memory management and reducing overhead.
- ❖ Improved concurrency: COW enables concurrent access to data while ensuring data consistency, as modifications are only made to private copies of the data.

MUHAMMAD BILAL RAFIQ

# PAPER: Operating Systems                    Course Code: IT-306

Q#36: Which system call is used to create a process in UNIX?

Ans: The system call used to create a process in UNIX is fork().

fork() creates a new process that is a copy of the calling process. The new process inherits the parent process's memory space, file descriptors, and other resources. After the fork() call, the parent process and the child process both continue execution concurrently.

**How fork() works:**

1. The kernel allocates a new process control block (PCB) for the child process.
2. The kernel copies the parent process's memory space to the child process's memory space.
3. The child process inherits the parent process's file descriptors and other resources.
4. The child process starts execution at the same instruction as the parent process.

**Return value of fork():**

❖ On success: The parent process returns the process ID (PID) of the child process.
❖ On error: The parent process returns -1, and the child process is not created.

Example:

```c
#include <stdio.h>
#include <unistd.h>

int main() {
    pid_t pid = fork();

    if (pid == 0) {
        // Child process
        printf("Child process with PID: %d\n", getpid());
    } else if (pid > 0) {
        // Parent process
        printf("Parent process with PID: %d\n", getpid());
    } else {
        // Error
        printf("Error creating child process\n");
    }

    return 0;
}
```

MUHAMMAD BILAL RAFIQ

**Q#37: How hardware devices use functionality of an operating system?**

**Ans**: Hardware devices and operating systems have a symbiotic relationship. While hardware provides the physical components for a computer system, the operating system acts as the intermediary between the hardware and software applications, enabling them to interact and function effectively. Here's how hardware devices utilize the functionalities of an operating system:

1. **Device Drivers:**
   - ❖ Operating systems provide device drivers, which are specialized software programs that allow the OS to communicate with specific hardware devices.
   - ❖ Examples: Graphics card drivers, network card drivers, printer drivers, etc.
2. **Interrupt Handling:**
   - ❖ Hardware devices can generate interrupts, which are signals sent to the OS to notify it of events requiring attention.
   - ❖ Examples: Keyboard input, mouse clicks, network data arrival, etc.
3. **Memory Management:**
   - ❖ The OS manages the allocation and deallocation of memory for both the device drivers and the data buffers used by the devices.
   - ❖ Examples: Allocating memory for a graphics card's frame buffer, managing data buffers for network communication, etc.
4. **I/O Management:**
   - ❖ The OS provides a unified interface for applications to perform input/output operations with various devices.
   - ❖ Examples: Reading data from a keyboard, writing data to a printer, sending data over a network, etc.
5. **Power Management:**
   - ❖ The OS can control the power state of hardware devices to optimize energy consumption.
   - ❖ Examples: Putting the hard disk to sleep after a period of inactivity, dimming the display when the computer is idle, etc.
6. **Device Security:**
   - ❖ The OS can implement security measures to protect devices from unauthorized access and malicious attacks.
   - ❖ Examples: Restricting access to specific devices based on user privileges, encrypting data transmitted over a network, isolating virtual machines from the host system's hardware, etc.

# PAPER: Operating Systems                    Course Code: IT-306

**Q#38: What is meant by CPU-protection?**

**Ans:** CPU protection refers to a set of mechanisms and techniques implemented in hardware and software to safeguard the central processing unit (CPU) from unauthorized access, malicious code execution, and other security threats. It encompasses various aspects of CPU security, including:

1. **Memory Protection:**
   - Prevents unauthorized processes from accessing or modifying memory regions that belong to other processes or the operating system.
   - Achieved through techniques like memory segmentation, paging, and access control mechanisms.

2. **Instruction Protection:**
   - Safeguards the CPU from executing malicious instructions or unauthorized code.
   - Implemented through mechanisms like code signing, instruction set extensions for security, and privilege levels.

3. **Data Protection:**
   - Protects sensitive data from unauthorized access or modification during processing within the CPU.
   - Achieved through techniques like encryption, data isolation, and secure memory management.

4. **Exception Handling:**
   - Provides a mechanism for the CPU to handle unexpected events or errors securely.
   - Prevents malicious code from exploiting exceptions to gain control or access sensitive data.

5. **Virtualization Support:**
   - Enables the creation of isolated virtual environments that share the same physical CPU resources.
   - Provides a layer of protection between different virtual machines and the host system.

**Examples of CPU Protection Mechanisms:**

- **Memory Management Units (MMUs):** Hardware units that enforce memory protection by translating virtual addresses to physical addresses and checking access permissions.
- **x86 Privilege Levels:** Different privilege levels for instructions and data, allowing the operating system to control access to sensitive resources.
- **ARM TrustZone:** A hardware-based security extension that creates isolated execution environments for sensitive code and data.
- **Intel SGX (Software Guard Extensions):** Provides enclaves for secure code execution and data protection.

MUHAMMAD BILAL RAFIQ

**Q#39: What is meant by dual-mode operation?**

**Ans:** Dual-mode operation is a fundamental concept in operating systems that refers to the ability of the CPU to operate in two distinct modes:

1. **Kernel Mode (Supervisor Mode):**
   - ❖ Privileged Mode: The CPU has unrestricted access to all system resources, including memory, I/O devices, and other privileged instructions.
   - ❖ Used for: Executing operating system code, handling interrupts, managing system resources, and performing other privileged tasks.
   - ❖ Security Risk: Code running in kernel mode can potentially compromise the entire system if not carefully designed and implemented.
2. **User Mode (Problem Mode):**
   - ❖ Non-Privileged Mode: The CPU has restricted access to system resources and cannot execute privileged instructions.
   - ❖ Used for: Executing user applications and programs.
   - ❖ Security Benefit: Limits the damage that can be caused by malicious or buggy user applications.

**Examples of Dual-Mode Operation:**

- ❖ x86 Architecture: Uses four privilege levels, with ring 0 being the most privileged (kernel mode) and ring 3 being the least privileged (user mode).
- ❖ ARM Architecture: Uses two privilege levels: privileged mode (kernel mode) and user mode.

**Q#40: What is the function of ready queue?**

**Ans:** The ready queue, also known as the CPU ready queue, is a data structure in a computer's operating system that holds all the processes that are currently in the ready state, waiting for their turn to be executed by the central processing unit (CPU). It acts as a kind of staging area or buffer between the processes that are running and the rest of the system, ensuring that processes are managed fairly and efficiently.

The scheduler, a part of the operating system,  uses the ready queue to determine which process is next in line for execution based on a scheduling algorithm that considers factors like priority, time, and other criteria.  This ensures that the CPU is being used efficiently and that no processes are left waiting indefinitely.

MUHAMMAD BILAL RAFIQ

# PAPER: Operating Systems                    Course Code: IT-306

**Q#41: In which memory management technique process can be divided in to equal parts?**

**Ans**: The memory management technique that divides the process into equal-sized blocks is called paging.

In paging, both the process address space and the physical memory are divided into fixed-size blocks called pages and frames, respectively. The size of a page is typically a power of 2, such as 4 KB or 8 KB. This simplifies memory management and reduces external fragmentation.

Here's a breakdown of how paging works:

- ❖ **Process division**: The process is divided into equal-sized pages.
- ❖ **Physical memory division**: The physical memory is divided into equal-sized frames.
- ❖ **Page table**: A page table is maintained that keeps track of the mapping between a process's pages and the available frames in physical memory.
- ❖ **Address translation**: When a process generates a memory address (logical address), the Memory Management Unit (MMU) translates it into a physical address using the page table.

**Q#42: In which memory management technique, system suffers from External Fragmentation?**

**Ans:** There are two main memory management techniques that suffer from external fragmentation:

**Contiguous Allocation with Variable-Sized Partitions:** This technique allocates memory to processes in contiguous blocks, but the size of the blocks can vary depending on the process's needs. While this approach allows for better memory utilization compared to fixed-sized partitions, it leads to external fragmentation over time. As processes are loaded and unloaded, free memory gets broken up into smaller chunks that may not be large enough to fit new processes, even though there might be sufficient total free space available.

**Segmentation (Pure):** Segmentation divides a process's logical address space into variable-sized segments based on its logical structure (code, data, stack etc.). This offers flexibility in memory allocation, but similar to variable-sized contiguous allocation, it can lead to external fragmentation. Free memory becomes scattered in chunks of various sizes that may not be suitable for new processes requiring memory.

**Q#42: What is meant by critical section?**

**Ans:** In an operating system, a critical section refers to a specific block of code within a program that deals with shared resources. These shared resources can be various things, like:

- ❖ Memory locations: This could be a variable or data structure that multiple processes or threads need to access and modify.
- ❖ Files: When multiple processes need to read from or write to the same file concurrently.
- ❖ Hardware devices: Imagine multiple processes trying to print to the same printer at the same time.

The critical part is that if multiple processes or threads try to access and modify these shared resources simultaneously, it can lead to inconsistencies and errors. This is called a race condition.

MUHAMMAD BILAL RAFIQ

**Q#43: What is meant by race condition?**

**Ans:** A race condition is a situation in concurrent programming where the outcome of a computation depends on the unpredictable timing of multiple processes or threads accessing shared resources. This can lead to inconsistent and erroneous results, as the order in which processes access and modify shared data becomes critical.

**Causes of Race Conditions:**

- ❖ **Shared resources:** Multiple processes or threads accessing and modifying the same shared data, such as global variables, data structures, or hardware devices.
- ❖ **Unpredictable timing:** The timing of process execution and resource access is unpredictable, leading to non-deterministic behavior.
- ❖ **Lack of synchronization:** Processes or threads do not synchronize their access to shared resources, resulting in potential conflicts and race conditions.

Q#44: Mention four reasons for process creation?

Ans: Four Reasons for Process Creation

1. **Concurrency**: To achieve concurrency, multiple processes can be created to execute different tasks simultaneously or concurrently. This allows for efficient utilization of system resources and can improve the responsiveness of an application.
2. **Modularity**: Processes can be created to represent different modules or components of a larger program. This enhances modularity and makes it easier to develop, debug, and maintain the program.
3. **Resource Sharing:** Processes can share resources, such as memory, files, and devices, through inter-process communication (IPC) mechanisms. This allows for efficient resource utilization and collaboration between processes.
4. **Protection and Isolation:** Processes can be created to provide protection and isolation between different parts of a program or system. This can help prevent errors in one process from affecting other processes or the operating system.

**Q#45: What is difference between Load Time Dynamic Linking VS Run Time Dynamic Linking?**

**Ans:** Dynamic linking is a technique that allows a program to load and link with libraries or modules at runtime instead of compile time. This provides several advantages, including reduced memory footprint, modularity, and easier updates. However, there are two main types of dynamic linking: load-time dynamic linking and run-time dynamic linking.

**Load-Time Dynamic Linking (LTLD):**

❖ **Binding**: Libraries or modules are linked with the program at load time, typically when the program starts.
❖ **Mechanism**: The linker resolves external references and creates a table of function pointers within the program's memory space.
❖ **Advantages**:
- Faster startup time compared to run-time linking.
- Reduced memory footprint as only the required libraries are loaded.
- Easier debugging as symbols are available at load time.
❖ **Disadvantages**:
- Requires static linking information about the libraries.
- Libraries cannot be dynamically loaded or unloaded at runtime.

**Run-Time Dynamic Linking (RTLD):**

❖ **Binding**: Libraries or modules are linked with the program at runtime, typically when a function from the library is first called.
❖ **Mechanism**: The program dynamically loads the library and resolves external references when needed.
❖ **Advantages:**
- Allows for dynamic loading and unloading of libraries at runtime.
- Enables lazy loading of libraries, reducing startup time.
- Provides greater flexibility for modularity and updates.
❖ **Disadvantages**:
- Slower startup time compared to load-time linking.
- Requires additional runtime overhead for dynamic linking.
- Debugging can be more challenging due to late binding.

**Examples**:

✓ **LTLD**: Windows DLLs, statically linked shared libraries in Linux.
✓ **RTLD**: Dynamically loaded shared libraries in Linux, Java class loading.

MUHAMMAD BILAL RAFIQ

# Long Questions

**Q#46: Recall reader-writer problem with reader's priority and check the following code. Re by adding missing lines.**

```
do {
        readcount++;
        if (readcount == 1)        wait(wrt);

        ...
        // reading is performed
        ...
        readcount--;
        if (readcount == 0)
                signal(wrt);
} while(1);
```

**Ans:** Here's a breakdown of the code and the missing lines:

```
1  int readcount = 0;  // Initialize reader count to 0
2  semaphore wrt = 1;   // Initialize semaphore wrt to 1 (unlocked)
3
4  do {
5     readcount++;
6     wait(wrt);
7
8     // Reader is in the critical section (your code to read goes here)
9
10    readcount--;
11    if (readcount == 0) {
12       signal(wrt);
13    }
14 } while(1);
15
```

Key points of reader priority:

- ❖ Readers can access the shared resource concurrently as long as no writer is accessing it.
- ❖ A writer has to wait for all readers to finish before it can enter the critical section.

This solution ensures that readers do not starve (wait indefinitely) but it can lead to writer starvation if there are many readers and few writers.

*MUHAMMAD BILAL RAFIQ*

# PAPER: Operating Systems                    Course Code: IT-306

**Q#47: Consider a logical address space of 64 pages of 1024 words each, mapped onto a physical memory of 32 frames.**

   a.  **How many bits are there in the logical address?**
   b.  **How many bits are there in the physical address?**
   c.  **How many bits needed for offset?**

**Ans**: Logical Address Space and Physical Memory Analysis

   a.  **Bits in the logical address:**
       ❖  We have 64 pages, each with 1024 words.
       ❖  Total words in logical address space = 64 pages * 1024 words/page = 65536 words.
       ❖  To address 65536 words, we need 2^16 bits (16 bits).
   b.  **Bits in the physical address:**
       ❖  We have 32 frames.
       ❖  To address 32 frames, we need 2^5 bits (5 bits).
   c.  **Bits needed for offset:**
       ❖  Each page has 1024 words.
       ❖  To address 1024 words within a page, we need 2^10 bits (10 bits).

**Therefore:**

       ❖  Logical address has 16 bits.
       ❖  Physical address has 5 bits.
       ❖  Offset needs 10 bits.

MUHAMMAD BILAL RAFIQ

**Q#48: A system has four process and five resources. The current allocation and maximum needs are as follows:**

|   | Allocated | | | | | Maximum | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| A | 1 | 0 | 2 | 1 | 1 | 1 | 1 | 2 | 1 | 3 |
| B | 2 | 0 | 1 | 1 | 0 | 2 | 2 | 2 | 1 | 0 |
| C | 1 | 1 | 0 | 1 | 1 | 2 | 1 | 3 | 1 | 1 |
| D | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 2 | 2 | 0 |

**If Available matrix consist of [0,0,2,1,1], check whether system is in safe state are not. Also write safe sequence if system is in safe state.**

**Ans**: First, we calculate the Need matrix.

Need = Max – Allocated

$$Need = \begin{matrix} 1\ 2\ 1\ 3 \\ 2\ 2\ 1\ 0 \\ 1\ 3\ 1\ 1 \\ 1\ 2\ 2\ 0 \end{matrix}$$

Then, we calculate the Available matrix.

Available = Total Resources – Allocated Resources

Available = 0 0 2 1 1

Now, we can check whether the system is in a safe state or not.

We will use Banker's algorithm to check for safety.

1. Initialize Finish[] = {false, false, false, false}
2. While(∃ i such that Finish[i] == false and Need[i] ≤ Available)
   a. Finish[i] = true
   b. Available = Available + Allocated[i]
3. If(∀ i, Finish[i] == true)
   a. System is in safe state
   b. Print Safe Sequence
4. Else
   a. System is not in safe state

Let's apply the algorithm to the given example.

*MUHAMMAD BILAL RAFIQ*

**Iteration 1:**

- ❖  i = 0
- ❖  Finish[0] = true
- ❖  Available = Available + Allocated[0]
- ❖  Available = [0+1, 0+0, 2+2, 1+1, 1+1]
- ❖  Available = [1, 0, 4, 2, 2]

**Iteration 2:**

- ❖  i = 1
- ❖  Finish[1] = true
- ❖  Available = Available + Allocated[1]
- ❖  Available = [1+2, 0+0, 4+1, 2+1, 2+0]
- ❖  Available = [3, 0, 5, 3, 2]

**Iteration 3:**

- ❖  i = 2
- ❖  Finish[2] = true
- ❖  Available = Available + Allocated[2]
- ❖  Available = [3+1, 0+1, 5+0, 3+1, 2+1]
- ❖  Available = [4, 1, 5, 4, 3]

**Iteration 4:**

- ❖  i = 3
- ❖  Finish[3] = true
- ❖  Available = Available + Allocated[3]
- ❖  Available = [4+1, 1+1, 5+1, 4+1, 3+0]
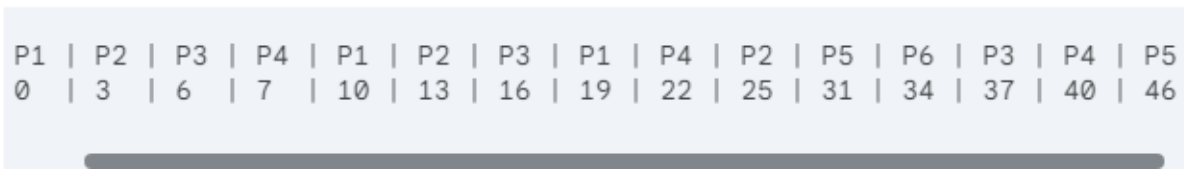- ❖  Available = [5, 2, 6, 5, 3]

Since all the processes are finished, the system is in a safe state.

**The safe sequence is <P0, P1, P2, P3>**

MUHAMMAD BILAL RAFIQ

**Q#49: Consider the following set of Processes, with the length of CPU time given in microseconds. Draw the Gantt chart, using Round Robin Algorithm when Time Slice = 3 microseconds. Also compute the waiting and turnaround times of each process.**

| Process | Arrival Time | Burst Time |
|---------|--------------|------------|
| P1 | 0 | 7 |
| P2 | 4 | 5 |
| P3 | 5 | 4 |
| P4 | 7 | 4 |
| P5 | 12 | 6 |
| P6 | 17 | 3 |

**Ans**: Here's the Gantt chart for the given set of processes using the Round Robin algorithm with a time slice of 3 microseconds:

```
P1 | P2 | P3 | P4 | P1 | P2 | P3 | P1 | P4 | P2 | P5 | P6 | P3 | P4 | P5
0  | 3  | 6  | 7  | 10 | 13 | 16 | 19 | 22 | 25 | 31 | 34 | 37 | 40 | 46
```

## Waiting Time:

- Process 1: (10 - 0) + (19 - 10) = 10 + 9 = 19 microseconds
- Process 2: (13 - 3) + (25 - 13) = 10 + 12 = 22 microseconds
- Process 3: (6 - 4) + (16 - 6) + (37 - 16) = 2 + 10 + 21 = 33 microseconds
- Process 4: (7 - 4) + (22 - 7) + (40 - 22) = 3 + 15 + 18 = 36 microseconds
- Process 5: (31 - 25) = 6 microseconds
- Process 6: (34 - 31) = 3 microseconds

## Turnaround Time:

- Process 1: 19 + 10 = 29 microseconds
- Process 2: 22 + 5 = 27 microseconds
- Process 3: 33 + 5 = 38 microseconds
- Process 4: 36 + 7 = 43 microseconds
- Process 5: 6 + 12 = 18 microseconds
- Process 6: 3 + 17 = 20 microseconds

MUHAMMAD BILAL RAFIQ

**Q#50: Consider three resource types in a system with following number of instances.**

**A (14 instances), B (10 instances), C (20 instances)**

**System State is as shown below, find safe sequence using Bankers algorithms.**

| PROCESS | MAX | | | ALLOCATION | | | AVAILABLE | | |
|---------|---|---|---|---|---|---|---|---|---|
| | A | B | C | A | B | C | A | B | C |
| P0 | 6 | 10 | 14 | 0 | 2 | 0 | 4 | 6 | 6 |
| P1 | 4 | 4 | 6 | 0 | 0 | 4 | | | |
| P2 | 4 | 0 | 18 | 4 | 0 | 6 | | | |
| P3 | 4 | 4 | 4 | 2 | 2 | 4 | | | |
| P4 | 6 | 6 | 8 | 4 | 0 | 0 | | | |

**Ans**: PROCESS      Safe sequence is <P1, P3, P4, P0, P2>

Here is the detail of each iteration:

Iteration 1:

Process P1 is requesting for 2 instances of resource A, 0 instances of resource B and 0 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

| apache | | | | | | | | | |
|--------|---|---|---|---|---|---|---|---|---|
| PROCESS | MAX | | | ALLOCATION | | | AVAILABLE | | |
| A B C | | A | B | C | | A | B | C | |
| P0 | 6 | 10 | 14 | 0 | 2 | 0 | 4 | 6 | 6 |
| P1 | 4 | 4 | 6 | 2 | 4 | 0 | 4 | 2 | 6 |
| P2 | 4 | 0 | 18 | 4 | 0 | 6 | 0 | 6 | 0 |
| P3 | 4 | 4 | 4 | 2 | 2 | 4 | 0 | 2 | 0 |
| P4 | 6 | 6 | 8 | 4 | 0 | 0 | 0 | 6 | 8 |

**Iteration 2:**

Process P3 is requesting for 2 instances of resource A, 2 instances of resource B and 0 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

| apache | | | | | | | | | |
|--------|---|---|---|---|---|---|---|---|---|
| PROCESS | MAX | | | ALLOCATION | | | AVAILABLE | | |
| A B C | | A | B | C | | A | B | C | |
| P0 | 6 | 10 | 14 | 0 | 2 | 0 | 2 | 4 | 6 |
| P1 | 4 | 4 | 6 | 2 | 4 | 0 | 2 | 0 | 6 |
| P2 | 4 | 0 | 18 | 4 | 0 | 6 | 0 | 6 | 0 |
| P3 | 4 | 4 | 4 | 4 | 4 | 4 | 0 | 0 | 0 |
| P4 | 6 | 6 | 8 | 4 | 0 | 0 | 0 | 6 | 8 |

*MUHAMMAD BILAL RAFIQ*

**Iteration 3:**

Process P4 is requesting for 4 instances of resource A, 0 instances of resource B and 0 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

apache

| PROCESS | | | MAX | | | ALLOCATION | | | AVAILABLE | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | B | C | A | B | C | A | B | C | | | |
| P0 | | 6 | 10 | 14 | | 0 | 2 | 0 | 6 | 4 | 6 |
| P1 | | 4 | 4 | 6 | | 2 | 4 | 0 | 6 | 0 | 6 |
| P2 | | 4 | 0 | 18 | | 4 | 0 | 6 | 6 | 6 | 0 |
| P3 | | 4 | 4 | 4 | | 4 | 4 | 4 | 6 | 0 | 0 |
| P4 | | 6 | 6 | 8 | | 8 | 0 | 0 | 2 | 6 | 8 |

**Iteration 4:**

Process P0 is requesting for 6 instances of resource A, 8 instances of resource B and 14 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

apache

| PROCESS | | | MAX | | | ALLOCATION | | | AVAILABLE | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| A | B | C | A | B | C | A | B | C | | | |
| P0 | | 6 | 10 | 14 | | 6 | 10 | 14 | 0 | 0 | 0 |
| P1 | | 4 | 4 | 6 | | 2 | 4 | 0 | 0 | 0 | 6 |
| P2 | | 4 | 0 | 18 | | 4 | 0 | 6 | 0 | 6 | 0 |
| P3 | | 4 | 4 | 4 | | 4 | 4 | 4 | 0 | 0 | 0 |
| P4 | | 6 | 6 | 8 | | 8 | 0 | 0 | 0 | 6 | 8 |

**Process P0 is now finished, so it is removed from the system. The available resources are updated as follows:**

apache

AVAILABLE = (4 + 6) + (6 + 8) + (6 + 14) = 16 + 14 + 20 = 50

MUHAMMAD BILAL RAFIQ

**Iteration 5:**

Process P2 is requesting for 4 instances of resource A, 0 instances of resource B and 18 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

```
apache

PROCESS         MAX         ALLOCATION          AVAILABLE
A   B   C       A   B   C       A   B   C         A   B   C
P1      4   4   6       2   4   0       16   6   8
P2      4   0   18      8   0   18       8   6   0
P3      4   4   4       4   4   4        8   0   0
P4      6   6   8       8   0   0        8   6   8
```

Process P2 is now finished, so it is removed from the system. The available resources are updated as follows:

```
apache                                              Copy

AVAILABLE = (8 + 4) + (6 + 0) + (18 + 0) = 12 + 6 + 18 = 36
```

Iteration 6:
Process P3 is requesting for 2 instances of resource A, 2 instances of resource B and 0 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

```
apache                                              Copy

PROCESS     MAX     ALLOCATION     AVAILABLE
A  B  C     A  B  C     A  B  C
P1      4  4  6     2  4  0     14  2  8
P3      4  4  4     6  6  4      8  0  0
P4      6  6  8     8  0  0      8  6  8
```

*MUHAMMAD BILAL RAFIQ*

Process P3 is now finished, so it is removed from the system. The available resources are updated as follows:

```apache
AVAILABLE = (6 + 2) + (4 + 2) + (4 + 0) = 8 + 6 + 4 = 18
```

Iteration 7:

Process P4 is requesting for 2 instances of resource A, 6 instances of resource B and 8 instances of resource C. Since all the instances are available, the request is granted. The allocation table is updated as follows:

```apache
PROCESS     MAX      ALLOCATION     AVAILABLE
A  B  C     A  B  C      A  B  C
P1     4  4  6     2  4  0     16  0  8
P4     6  6  8    10  6  8      6  0  0
```

Process P4 is now finished, so it is removed from the system. The available resources are updated as follows:

```apache
AVAILABLE = (10 + 6) + (6 + 0) + (8 + 0) = 16 + 6 + 8 = 30
```

Since all the processes are finished, the system is in a safe state. The safe sequence is <P1, P3, P4, P0, P2>.

*MUHAMMAD BILAL RAFIQ*

**Q#51: Consider a system with 64-bit logical address. This system implements paging with 2 KB page size. Compute: lengths of p and d fields in the logical address.**

**Ans**: Compute the lengths of the page number (p) and offset (d) fields in the logical address for a system with 64-bit logical address and 2 KB page size:

1. **Page Size**: We are given that the page size is 2 KB. In bits, this translates to:
   - Page size (bits) = 2 KB * 8 bits/byte = 2 * 1024 bits = $2^{11}$ bits
2. **Logical Address Size**: The logical address is 64 bits.
3. **Relationship between p and d**: The logical address is divided into two parts: the page number (p) and the offset (d) within the page.
4. **Computing p**: Since the offset specifies the address within a page of size $2^{11}$ bits, the remaining bits in the logical address must represent the page number. So, we can find the length of p as follows:
   - Logical Address size (bits) - Offset size (bits) = p (bits)
   - 64 bits - $2^{11}$ bits = p (bits)
5. **Calculating p:**
   - p (bits) = 64 bits - 2048 bits = 53 bits
   - Therefore, the page number field (p) in the logical address has a length of 53 bits.
6. **Offset (d):** As we already established, the offset (d) specifies the address within a page and its size is equal to the page size in bits:
   - d (bits) = Page size (bits)
   - d (bits) = $2^{11}$ bits = 2048 bits

**Q#52: In a paging system, logical address space is 4 GB and the page size is 4 KB. Available physical memory is 16 GB. Compute: length of logical address and length of f field in the physical address.**

**Ans:** Compute the lengths of the logical address and the frame number (f) field in the physical address for a paging system with the given specifications:

1. **Logical Address Space**: We are given that the logical address space is 4 GB.
2. **Logical Address Size**: To find the size of the logical address in bits, we need to convert gigabytes (GB) to bits:
   - Logical Address size (bits) = 4 GB * (8 bits/byte) * (1024 bytes/KB) * (1024 KB/MB)
   - Note: We use the conversion factors of 8 bits/byte and 1024 units (bytes or kilobytes) per unit (kilobyte or megabyte) to convert between units.
3. **Calculation**:
   - Logical Address size (bits) = 4 * 8 * 1024 * 1024 = $2^{32}$ bits
   - Page Size: The page size is given as 4 KB. Converting this to bits:
   - Page size (bits) = 4 KB * 8 bits/byte = $2^{12}$ bits
4. **Relationship between Logical Address and Frame Number**: In paging, the logical address is divided into two parts: the page number (p) and the offset (d) within the page. The frame number (f) in the physical address corresponds to the page number in the logical address.

MUHAMMAD BILAL RAFIQ

5. **Computing f:** Since the offset specifies the address within a page of size 2^12 bits, the remaining bits in the logical address must represent the frame number. So, we can find the length of f as follows:
   - ❖ Logical Address size (bits) - Offset size (bits) = f (bits)
   - ❖ 2^32 bits - 2^12 bits = f (bits)
6. **Calculating f:**
   - ❖ f (bits) = 2^32 bits - 2^12 bits = 2^20 bits

**Q#53: If the hit-ratio to a TLB is 98%, and it takes 10 nanoseconds (nsec) to search the TLB and 100 nsec to access the main memory, then what must be the Effective Memory Access Time in nanoseconds?**

**Ans**: Calculate the Effective Memory Access Time (EMAT) for the system:

1. **Hit Ratio**: We are given that the hit ratio to the TLB (Translation Lookaside Buffer) is 98%. This means in 98% of cases, the memory address translation happens using the TLB, which is faster than main memory access.
2. **Miss Ratio**: The remaining 2% of cases represent misses in the TLB, requiring access to the main memory. So, the miss ratio is 100% - 98% = 2%.
3. **Access Times**:
   - ❖ TLB access time: 10 nanoseconds (nsec)
   - ❖ Main memory access time: 100 nanoseconds (nsec)
4. **Weighted Average**: The EMAT is essentially a weighted average of the access times based on the probability of each scenario (hit or miss).
5. **Calculation**:
   - ❖ EMAT = (Hit ratio * TLB access time) + (Miss ratio * Main memory access time)
6. **Plugging in the values**:
   - ❖ EMAT = (0.98 * 10 nsec) + (0.02 * 100 nsec)
7. **Calculation:**
   - ❖ EMAT = 9.8 nsec + 2 nsec = 11.8 nsec (approximately)

**Therefore, the Effective Memory Access Time (EMAT) for this system is approximately 11.8 nanoseconds.**

**Q#54: If the total number of available frame is 50, and there are 2 processes; one of 50 pages and the other of 35 pages, then how much memory would be proportionally allocated to each of these processes?**

**Ans:** We can calculate the proportional memory allocation for each process using the following steps:

1. **Total Available Frames**: We are given that there are 50 total available frames in the system.
2. **Process Sizes**: There are two processes:
   - ❖ Process 1: 50 pages
   - ❖ Process 2: 35 pages
3. **Proportional Allocation**: We want to allocate memory based on the ratio of each process's size to the total process size.
4. **Total Process Size:**
   - ❖ Total process size = Size of process 1 + Size of process 2
   - ❖ = 50 pages + 35 pages
   - ❖ = 85 pages
5. **Proportional Allocation Ratio**:
   - ❖ For process 1: (Size of process 1) / (Total process size)
   - ❖ = 50 pages / 85 pages
   - ❖ For process 2: (Size of process 2) / (Total process size)
   - ❖ = 35 pages / 85 pages

**Since these ratios are in fractions, it might be easier to work with them in decimal form.**

6. **Proportional Allocation (Decimal):**
   - ❖ Process 1: 50 pages / 85 pages ≈ 0.5882
   - ❖ Process 2: 35 pages / 85 pages ≈ 0.4118
7. **Memory Allocation (Frames):** In a perfect scenario, we would allocate these ratios directly to the number of frames. However, since the number of frames needs to be whole numbers, we might need to round the results.
   - ❖ Process 1 (frames): 0.5882 * 50 frames ≈ 29.41 (rounded up to 29 frames)
   - ❖ Process 2 (frames): 0.4118 * 50 frames ≈ 20.59 (rounded up to 21 frames)

## Summary:

- Process 1: Proportionally allocated ~ 58.82% (rounded up to 29 frames)
- Process 2: Proportionally allocated ~ 41.18% (rounded up to 21 frames)

MUHAMMAD BILAL RAFIQ

**Q#55: What is meant by Belady's Anomaly? Why stack based algorithms do not suffer from Belady's Anomaly?**

**Ans:** Belady's Anomaly is a phenomenon in memory management that occurs in paging systems using certain page replacement algorithms. It refers to the counterintuitive situation where increasing the number of available page frames (memory) can actually lead to an increase in the number of page faults (accesses to the disk to retrieve data).

**Here's a breakdown of Belady's Anomaly:**

1. **Page Faults:** When a program tries to access data that is not currently loaded in physical memory, a page fault occurs. The operating system needs to fetch the required page from disk and load it into a free frame.
2. **Page Replacement Algorithms:** To manage memory efficiently, operating systems use page replacement algorithms. These algorithms decide which page frame to evict from memory when a new page needs to be loaded and there are no free frames available.
3. **Belady's Anomaly:** This anomaly occurs with specific page replacement algorithms, most notably the First-In-First-Out (FIFO) algorithm. In FIFO, the oldest page in memory (the one loaded first) is evicted to make space for a new page.

**Why FIFO Suffers from Belady's Anomaly:**

❖ FIFO doesn't consider the future usage of a page. It simply evicts the oldest page regardless of whether it will be needed soon.
❖ Belady's Anomaly can occur with FIFO when a sequence of memory accesses repeatedly uses a specific set of pages that barely fit in the available frames.
❖ Increasing the number of frames might initially seem beneficial, allowing more pages to be loaded. However, it can also lead to the eviction of recently used pages that will be needed again soon, causing more page faults.

**Why Stack-Based Algorithms Avoid Belady's Anomaly:**

Stack-based algorithms prioritize pages based on their recent usage.  Here are some examples:

❖ Least Recently Used (LRU): This algorithm keeps track of which pages are accessed most recently and evicts the least recently used one. Pages that are actively used are less likely to be evicted, even if they are older.
❖ Second Chance Algorithm: This is a variant of FIFO that gives a second chance to a page before evicting it. When a page fault occurs and all frames are full, the second chance algorithm checks if the oldest page has been modified. If not, it evicts the page. Otherwise, it sets a "second chance" bit and moves on to the next page in the FIFO order.

By prioritizing recently used pages, stack-based algorithms are less likely to evict pages that will be needed again soon, even if the number of frames increases. This helps to avoid Belady's Anomaly.

MUHAMMAD BILAL RAFIQ