

VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



REPORT ASSIGNMENT

OPERATING SYSTEM

Class: CC01 - HK251
INSTRUCTOR: NGUYEN PHUONG DUY

TOPIC: DESIGN SIMPLE OPERATING SYSTEM

| student's name | Student's ID |
|------------------|--------------|
| Le Phuc Khang | 235 |
| Tran Quoc Thang | 2353125 |
| Tran Xuan Hao | 235 |
| Nguyen Xuan Ngoc | 235 |
| Phan Tuan Kiet | 2352654 |

GROUP WORK RESULTS REPORT

| Number | Full name | Student ID | Contribution |
|--------|------------------|------------|--------------|
| 1 | Le Phuc Khang | 235 | 100% |
| 2 | Tran Quoc Thang | 2353125 | 100% |
| 3 | Nguyen Xuan Ngoc | 235 | 100% |
| 4 | Phan Tuan Kiet | 2352654 | 100% |
| 5 | Tran Xuan Hao | 235 | 100% |

Contents

| | | |
|----------|---|-----------|
| 1 | Theoretical Background | 3 |
| 1.1 | Scheduler | 3 |
| 1.2 | Memory Management | 4 |
| 1.2.1 | Memory management | 4 |
| 1.2.2 | Multi-level paging | 5 |
| 1.2.3 | Comparative Analysis of Page Table Architectures | 6 |
| 2 | Implementation | 9 |
| 2.1 | Scheduler | 9 |
| 2.1.1 | Queue Operations (Implementation of Round-Robin) | 9 |
| 2.1.2 | MLQ Scheduler Implementation (sched.c) | 10 |
| 2.1.3 | Scheduler Algorithm Conclusion | 13 |
| 2.2 | Memory Management | 13 |
| 2.2.1 | Memory management | 13 |
| 2.2.2 | Multi-level paging | 18 |
| 3 | Interprets the results of running tests | 22 |
| 3.1 | Scheduling | 22 |
| 3.1.1 | The result input file of scheduler (sched) | 22 |
| 3.1.2 | Scheduling: Gantt chart | 24 |
| 3.2 | Memory | 24 |
| 3.2.1 | Result of the input file of Memory (os_0_mlq_paging) | 24 |
| 3.2.2 | Explain the status of the memory allocation in heap and data segments | 27 |
| 3.2.3 | Chronological Analysis of Memory Allocation | 28 |
| 3.2.4 | Summary of Result Flow | 29 |
| 4 | Answer question | 30 |
| 4.1 | Question 1 | 30 |
| 4.2 | Question 2 | 30 |
| 4.3 | Question 3 | 31 |
| 4.4 | Question 4 | 34 |
| 4.5 | Question 5 | 36 |
| 4.6 | Question 6 | 45 |
| 4.7 | Question 7 | 58 |
| 5 | Overall | 70 |
| 6 | Source Code | 70 |

1 Theoretical Background

1.1 Scheduler

The scheduler is the central component of the kernel responsible for deciding which process runs at any given time. Our implementation is grounded in the theoretical framework of CPU Scheduling described in Chapter 5 of *Operating System Concepts*. Specifically, we construct a **Symmetric Multiprocessing (SMP)** system utilizing a **Preemptive Multilevel Queue (MLQ)** algorithm with a **Time-Slicing** allocation strategy to address starvation.

1. Multilevel Queue Scheduling (Theory of Partitioning):

- **Theoretical Basis:** Processes in a system often have different response-time requirements (e.g., interactive/foreground vs. batch/background). A Multilevel Queue (MLQ) algorithm partitions the Ready Queue into several separate queues, each assigned to a different class of processes.
- **System Design:** In our assignment, we define `MAX_PRIO` distinct queues. Processes are permanently assigned to a specific queue based on their `priority` attribute. Unlike Multilevel Feedback Queues (MLFQ), which allow processes to move between queues to define their behavior dynamically, our standard MLQ implementation reduces runtime overhead ($\mathcal{O}(1)$ queue selection) by maintaining fixed process classification.

2. Intra-Queue Scheduling: Round-Robin (RR) and Preemption:

- **Theoretical Basis:** To support time-sharing, Round-Robin (RR) scheduling is theoretically defined as FCFS scheduling with **Preemption**. A small unit of time, called a *Time Quantum*, is defined. If a process's CPU burst exceeds this quantum, the system timer generates an interrupt, causing the OS to preempt the process and move it to the tail of the ready queue.
- **System Design:** We implement this by assigning a `time_slot` to the system. The `queue.c` module strictly enforces FIFO behavior (`dequeue` from head, `enqueue` to tail), which, combined with the timer interrupt in `os.c`, realizes the Round-Robin logic. This ensures **fairness** among processes of the same priority and minimizes the average response time.

3. Inter-Queue Scheduling: Solving Starvation via Time-Slicing:

- **The Starvation Problem:** A fundamental flaw of fixed-priority scheduling is **Indefinite Blocking** (or Starvation), where low-priority tasks may never execute if a continuous stream of high-priority tasks arrives.
- **Theoretical Solution:** The textbook proposes two solutions: *Aging* or *Time Slicing among queues*. In the Time Slicing approach, each queue gets a certain portion of the CPU time (e.g., 80% for foreground, 20% for background).
- **System Design (Slot Mechanism):** We implement the **Time Slicing** solution using a deterministic slot mechanism. Each queue i is allocated a specific budget: $slot[i] = MAX_PRIO - i$. The scheduler iterates through queues; if a high-priority queue consumes its allocated slots, the scheduler is theoretically forced to "yield" to lower-priority queues, even if the high-priority queue is non-empty.

This converts the algorithm from "Absolute Priority" to "Proportional Share," mathematically guaranteeing that even the lowest priority queue eventually receives CPU cycles.

4. Symmetric Multiprocessing (SMP) and Concurrency Control:

- **Theoretical Basis:** According to the theory of *Multiple-Processor Scheduling*, modern systems predominantly use Symmetric Multiprocessing (SMP), where each processor is self-scheduling. The literature outlines two primary strategies for organizing threads: (1) Per-processor private queues, and (2) A single **Common Ready Queue**. While private queues offer cache efficiency, the Common Ready Queue approach provides a unified pool of work but introduces specific synchronization challenges.
- **System Design (Load Balancing):** We adopt the **Common Ready Queue** architecture. As noted in the theoretical framework, this design naturally handles the distribution of workloads. Since all processors examine the same global queue to select a thread, the system inherently avoids the workload disparity often found in private-queue architectures. No complex balancing algorithms are needed because no processor remains idle while tasks exist in the global pool.
- **Synchronization Theory:** The textbook explicitly warns that a shared ready queue creates a potential **Race Condition**, where two separate processors might attempt to schedule the same thread simultaneously or corrupt the queue structure. To mitigate this, as mandated by the principles of concurrency control, we implement a locking mechanism (`pthread_mutex_t`). This lock protects the critical section (the shared queue), ensuring that enqueue and dequeue operations are strictly serialized, albeit at the cost of potential lock contention.

1.2 Memory Management

1.2.1 Memory management

The fundamental task of an operating system is to manage the memory hierarchy and provide an abstraction mechanism to separate the user process from the physical details of the hardware.

a. Logical vs. Physical Address Space According to **Chapter 9.1.3**, the central concept of memory management is the separation of logical address space from physical address space:

- **Logical Address:** An address generated by the CPU during program execution. To a process, memory appears as a flat, continuous space ranging from address 0 to $Max_{Virtual}$. The set of all valid logical addresses for a process is defined as the *Logical Address Space*.
- **Physical Address:** The actual address seen and accessed by the memory unit on the physical RAM.
- **Memory Management Unit (MMU):** A hardware device responsible for the run-time mapping from virtual to physical addresses. Under this scheme, the user program never interacts with the real physical addresses.

b. Process Memory Layout Although the virtual address space is linear, logically a process is typically organized into distinct sections to separate access rights and usage purposes (Chapter 9.3):

- **Text section:** contains the executable code of the program (usually read-only).
- **Data section:** stores global and static variables.
- **Heap:** region for dynamic allocations that conceptually grows upward.
- **Stack:** holds stack frames and local variables, growing downward.

In full-featured UNIX-like systems, each of the above logical sections is usually represented by one or more Virtual Memory Areas (VMAs) maintained in a linked list inside the process's `mm_struct`.

In our simple OS implementation, we adopt a slightly different but simpler model. Each process owns a single user `vm_area_struct` that covers its whole logical user space (`vma_id = 0`). Inside this VMA, we represent individual logical “segments” by region descriptors `vm_rg_struct`, whose start and end addresses (`rg_start`, `rg_end`) are recorded in the symbol region table `symrgtbl[]`. Free holes between regions are maintained in the VMA's `vm_freerg_list`.

Therefore, in this project:

- segment \approx `vm_rg_struct` region inside one VMA,
- while the VMA itself simply defines the outer bounds of the user address space for the process.

1.2.2 Multi-level paging

To implement memory mapping in modern architectures, **Paging** is widely used to eliminate external fragmentation.

a. The Page Table Problem in 64-bit Systems In traditional 32-bit architectures, a two-level paging scheme is often sufficient. However, with a 64-bit architecture (2^{64} bytes of address space), the size of the page table becomes prohibitively large if stored contiguously.

According to **Chapter 9.4.2 (Hierarchical Paging)**, the optimal solution is to divide the page table into smaller pieces. This allows the operating system to avoid allocating the entire page table structure in contiguous memory. Only the necessary sub-tables are allocated (on-demand allocation).

b. 5-Level Paging Architecture The system is designed based on an extension of the Intel x86-64 architecture, utilizing a **5-level paging scheme** to support a significantly larger virtual address space (up to 57-bit linear address).

The virtual address structure is decomposed into index fields and an offset. Assuming a standard page size of 4KB (2^{12} bytes), the virtual address is resolved as follows:

$$\text{Virtual_Address} = \text{PGD} \oplus \text{P4D} \oplus \text{PUD} \oplus \text{PMD} \oplus \text{PT} \oplus \text{Offset} \quad (1)$$

The address translation process (also known as the **Page Walk**) occurs sequentially through 5 tables:

1. **PGD (Page Global Directory):** The top-level table (Level 5).

2. **P4D (Page Level 4 Directory):** The Level 4 table.
3. **PUD (Page Upper Directory):** The Level 3 table.
4. **PMD (Page Middle Directory):** The Level 2 table.
5. **PT (Page Table):** The final level (Level 1), containing the mapping to the physical frame.

Each entry in the upper-level tables points to the base address of the next lower-level table. The entry in the PT (Page Table Entry - PTE) contains the actual **Physical Frame Number (PFN)**.

1.2.3 Comparative Analysis of Page Table Architectures

a. Hierarchical Paging This method divides the logical address into multiple parts to manage the page table as a tree structure (4-level or 5-level paging in modern systems).

Advantages

- **Memory Efficiency for Sparse Addresses:** The system does not need to allocate memory for the entire page table at once. It solves the problem where a single-level table would require an enormous amount of contiguous memory. It allows dividing the page table into smaller units.
- **Hardware Support:** This is the standard implementation for modern architectures, such as the Intel x86-64 and ARMv8 architectures.

Disadvantages

- **Memory Access Latency:** As the address space expands (e.g., to 64-bit), the number of levels in the hierarchy increases. Without a TLB, a single memory reference might require multiple memory accesses (e.g., 5 accesses) to traverse the page table before accessing the actual data.
- **Complexity:** Managing a multi-level tree (5 levels) is more complex than a simple flat table.

b. Hashed Page Tables This method uses a hash function on the virtual page number to locate the physical frame.

Advantages

- **Suitability for Large Address Spaces:** This structure is common in address spaces larger than 32 bits.
- **Handling Sparse Memory:** It is especially useful for sparse address spaces where memory references are non-contiguous and scattered.

Disadvantages

- **Collision Handling:** Multiple virtual addresses may hash to the same location, requiring a chain of elements to be searched to find a match.
- **Implementation Complexity:** While efficient for size, the overhead of hashing and traversing linked lists (chains) can impact performance compared to direct indexing.

c. Inverted Page Tables Instead of one page table per process, there is one entry for each real page (frame) of physical memory.

Advantages

- **Memory Savings:** This approach significantly decreases the amount of memory needed to store each page table, as the table size is proportional to physical memory, not the virtual address space.

Disadvantages

- **Increased Lookup Time:** It increases the time needed to search the table when a page reference occurs, often requiring a hash table to limit the search.
- **Difficulty with Shared Memory:** Implementing shared memory is difficult because there is typically only one mapping of a virtual address to a shared physical address.

d. Reasons for Choosing Hierarchical Page Tables In this simple operating system simulation project, we selected Hierarchical Page Tables (5-level Paging Scheme) for the following three reasons:

- **Dynamic Allocation:**

The hierarchical structure enables “on-demand” memory allocation. There is no need to initialize the entire massive data structure upfront. When a process requires access to a memory region, we simply invoke `alloc_aligned_table()` for that specific branch. This approach effectively prevents the wastage of management memory within a 64-bit environment.

- **Simplified Implementation of Sharing and Protection:**

The Address-Translation model allows for the assignment of access permissions (Read/Write/Execute) at each level of the tree (from top-level directories down to individual pages). This facilitates the implementation of Memory Protection, allowing for easy setup and granular control over specific memory regions.

- **Ease of Implementation:**

Compared to Inverted Page Tables or Hashed Page Tables, the hierarchical structure is easier to implement within a C-based simulation environment. We primarily rely on bitwise shifts and array indexing to decompose virtual addresses into their respective indices .

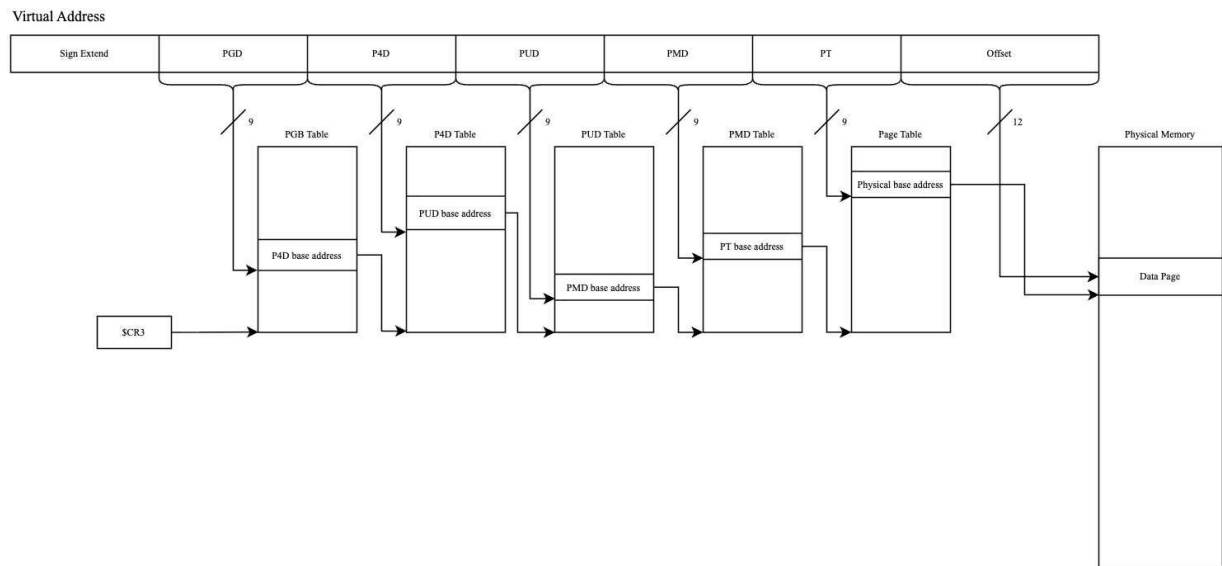


Figure 1: 5-Level Paging Scheme (Address Translation for 57-bit Virtual Address)

2 Implementation

2.1 Scheduler

2.1.1 Queue Operations (Implementation of Round-Robin)

To support the Round-Robin scheduling within each priority level, the queue operations in `queue.c` must strictly enforce First-In-First-Out (FIFO) behavior. I have implemented three core functions as follows:

a. Enqueue Operation (Arrival/Preemption) The `enqueue` function handles the arrival of a new process or the return of a preempted process.

- **Mechanism:** It places the process at index `q->size`. This is critical for Round-Robin: a process finishing its time slice is moved to the back of the line. We check for buffer overflow (`MAX_QUEUE_SIZE`) before assignment to ensure memory safety.

```
1 void enqueue(struct queue_t *q, struct pcb_t *proc) {
2     /* TODO: put a new process to queue [q] */
3     if (q->size < MAX_QUEUE_SIZE) {
4         q->proc[q->size] = proc;
5         q->size++;
6     }
7 }
```

Listing 1: Enqueue Implementation

b. Dequeue Operation (Dispatching) The `dequeue` function selects the next process for execution.

- **Mechanism:** It selects the process at index 0 (the Head), which corresponds to the process that has been waiting the longest. Since we use a static array, removing index 0 leaves a "hole". The function includes a loop to shift all subsequent elements ($i + 1$) to position i . This ensures the array remains contiguous and the next candidate is always at index 0.

```
1 struct pcb_t *dequeue(struct queue_t *q) {
2     /* TODO: return a pcb whose priority is the highest
3      * in the queue [q] and remember to remove it from q */
4     if (q == NULL || q->size == 0) {
5         return NULL;
6     }
7     struct pcb_t *first_proc = q->proc[0];
8
9     int i;
10    for (i = 0; i < q->size - 1; i++) {
11        q->proc[i] = q->proc[i + 1];
12    }
13    q->size--;
```

```
14     return first_proc;  
15 }
```

Listing 2: Dequeue Implementation

2.1.2 MLQ Scheduler Implementation (sched.c)

The `sched.c` module implements the core logic of the Multi-Level Queue scheduler. This implementation is designed to handle process retrieval and placement within a Symmetric Multiprocessing (SMP) environment, strictly following the Slot-based policy defined in the assignment specifications.

a. Scheduler Initialization (`init_scheduler`) This function sets up the system state before any process execution begins.

- **Slot Pre-calculation:** Instead of calculating the slot quota at runtime, we initialize the `slot` array immediately using the formula $slot[i] = MAX_PRIO - i$.
- **Reasoning:** This optimization reduces CPU overhead during the critical context-switching phase. By pre-filling the quotas, the scheduler simply needs to decrement or reset values during execution.
- **Mutex Initialization:** The `pthread_mutex_init` call prepares the lock required for thread safety in the SMP architecture.

```
1 void init_scheduler(void) {  
2     int i;  
3     for (i = 0; i < MAX_PRIO; i++) {  
4         mlq_ready_queue[i].size = 0;  
5         slot[i] = MAX_PRIO - i; // Pre-calculate quotas  
6     }  
7     // ...  
8     pthread_mutex_init(&queue_lock, NULL);  
9 }
```

Listing 3: Initialization Logic

b. Process Retrieval Logic (`get_mlq_proc`): The `get_mlq_proc` function serves as the core of the scheduler, responsible for selecting the next process for CPU allocation from the set of priority queues. The algorithm is designed based on the Multi-Level Queue (MLQ) model combined with a *Slot-based Allocation* mechanism to address the issue of resource starvation.

1. General Operating Principle

The function iterates sequentially through the queues from the highest priority (`prio = 0`) to the lowest (`prio = MAX_PRIO - 1`). The decision to select a process is based on two prerequisites:

- **Readiness:** The queue at that priority level must have waiting processes (`!empty`).
- **Resource Quota (Slot):** The queue must have remaining CPU usage quota (`slot[i] > 0`).

2. Algorithm Flow

The process selection procedure occurs through the following steps:

- **Step 1: Critical Section Protection.** Since the system simulates a multi-processor environment, the function begins by locking the Mutex (`pthread_mutex_lock(&queue_lock)`). This ensures the data integrity of the `mlq_ready_queue` when multiple CPUs access it simultaneously.
- **Step 2: Priority Traversal.** A loop iterates from $i = 0$ (highest priority) to `MAX_PRIO`. At each priority level i :
 - *Check Empty:* If `mlq_ready_queue[i]` is empty, the system skips it and checks the next priority level.
 - *Check Slot (Core Logic):*
 - * **Case A (Slot Remaining - $\text{slot}[i] > 0$):** The queue is allowed to run. The system retrieves the first process (`dequeue`) and decrements the remaining slots by 1 (`slot[i]--`).
Decision: Break the loop immediately to return this process to the CPU.
 - * **Case B (Slot Exhausted - $\text{slot}[i] == 0$):** The queue has used up its quota for the current cycle. The system performs a Slot Refill using the formula: `slot[i] = MAX_PRIO - i`.
Decision: Crucially, the algorithm DOES NOT select a process from this queue but continues the loop (`continue`) to check the next lower priority queue.
- **Step 3: System State Update.** If a process is found (`proc != NULL`), it is added to the running list (`running_list`) for tracking.
- **Step 4: Finish.** Unlock the Mutex (`pthread_mutex_unlock`) and return the `pcb_t` pointer of the selected process (or `NULL` if no viable process exists).

3. Slot-based Mechanism Analysis

This implements a mechanism of "**Conditional Yielding**":

- **Purpose:** To prevent high-priority processes (e.g., Priority 0) from monopolizing the CPU indefinitely.
- **Operation:** When a high-priority queue consumes its allocated slots (e.g., 140 slots for Prio 0), it is forced to reset its slots and "yield" the checking turn to lower-priority queues within the current function call. This ensures relative fairness, allowing lower-priority processes (such as Prio 139) a chance to execute (CPU time) even under heavy load conditions.

4. Complexity Analysis

In the worst-case scenario, the algorithm must traverse through all `MAX_PRIO` queues.

Complexity: $O(K + N)$

- K : The number of priority levels (`MAX_PRIO`).
- N : The cost of the `dequeue` operation.

```
1 struct pcb_t * get_mlq_proc(void) {
2     struct pcb_t * proc = NULL;
3
4     pthread_mutex_lock(&queue_lock);
5
6     for (int i = 0; i < MAX_PRIO; i++) {
7         if (!empty(&mlq_ready_queue[i])) {
8             if (slot[i] > 0) {
9                 proc = dequeue(&mlq_ready_queue[i]);
10                slot[i]--;
11                break;
12            } else {
13                slot[i] = MAX_PRIO - i;
14            }
15        }
16    }
17
18    if (proc != NULL)
19        enqueue(&running_list, proc);
20
21    pthread_mutex_unlock(&queue_lock);
22
23    return proc;
24 }
25 }
```

Listing 4: Slot-based Decision Logic

c. Process Placement (`put_mlq_proc` / `add_mlq_proc`) These functions handle putting a process back into the ready queue (e.g., after a time slice expires or upon creation).

- **Mechanism:** They map the process's priority (`proc->prio`) to the correct index in the `mlq_ready_queue` array and invoke the FIFO `enqueue` operation.
- **Synchronization:** These operations are wrapped in mutex locks to ensure the queue structure is not corrupted by concurrent access from the loader or other CPUs.

d. Synchronization Mechanism (Mutex Protection) A critical aspect of this implementation is the handling of Shared Resources in a multi-core environment.

In this simulation, the `mlq_ready_queue` is a **Global Shared Resource**.

- *Scenario without Mutex:* CPU 0 checks queue 0 and sees 1 process. Simultaneously, CPU 1 checks queue 0 and sees the same process. Both CPUs try to `dequeue` the same pointer. This creates a **Race Condition**, leading to memory corruption or one process being executed twice.

We use `pthread_mutex_lock(&queue_lock)` to define a **Critical Section**.

- **Scope:** The lock covers the entire duration of reading the state (checking `empty` and `slot`) and modifying the state (`dequeue`, `enqueue`, `slot --`).
- **Atomicity:** This ensures that the complex operation of "Check Slot \rightarrow Dequeue \rightarrow Decrement Slot" appears as a single atomic instruction to other CPUs. No other thread can interrupt or modify the queue during this sequence.

```
1 pthread_mutex_lock(&queue_lock); // BEGIN Critical Section
2 // ...
3 // Safe access to shared mlq_ready_queue and slot arrays
4 // ...
5 pthread_mutex_unlock(&queue_lock); // END Critical Section
```

Listing 5: Critical Section Protection

2.1.3 Scheduler Algorithm Conclusion

This section describes the implementation of a **Slot-based Multi-Level Queue (MLQ)** scheduler for an SMP environment. By integrating a quota mechanism ($slot = MAX_PRIO - prio$) into the priority traversal logic, the design **aims to balance** two conflicting objectives: prioritizing critical tasks while mitigating the starvation of lower-priority processes. Additionally, the application of **Mutex locks is designed to maintain** thread safety, **addressing** potential race conditions as multiple CPUs concurrently access the shared ready queue structure.

2.2 Memory Management

2.2.1 Memory management

a. Overview This module implements logical memory management per process in a simplified way. Each process owns an `mm_struct` that contains:

- One user `vm_area_struct` (`vmaid = 0`) acting as the container for the whole user address space,
- A fixed-size symbol region table `symrgtbl[]` that stores allocated regions by `rgid`,
- And the root pointers to the multi-level page tables handled by `mm64.c`.

Instead of multiple VMAs for text/data/heap/stack, our design keeps a single user VMA and represents logical “segments” as virtual memory regions (`vm_rg_struct`, VMR) inside that VMA. Each VMR is a contiguous range `[rg_start, rg_end)`, recorded in `symrgtbl[rgid]`, while free gaps are chained in `vm_freerg_list` for reuse.

The core design principles are:

- **VMA / region management:** All user regions live inside one user `vm_area_struct`; logical segments are `vm_rg_struct` descriptors, and free holes are kept in `vm_freerg_list`.

- **User / Kernel separation:** User code never touches kernel structures (`pcb_t`, `mm_struct`); it uses system calls and PID, and the kernel side looks up the correct PCB/MM before operating.
- **Allocation strategy:** Always try to reuse holes from `vm_freerg_list` first; if none fits, expand the VMA (`sbrk / vm_end`) via the `SYSMEM_INC_OP` syscall, then let the paging module map the new pages.

b. Implementation Details The following code snippets illustrate the core logic handling VMA organization and the secure User-Kernel interface.

1. VMA Data Structure (`os_mm.h`)

The `vm_area_struct` represents a contiguous range of logical addresses [`vm_start`, `vm_end`].

```
1 struct vm_area_struct {
2     unsigned long vm_id;      // Region ID (e.g., DATA, HEAP)
3     unsigned long vm_start;   // Logical Start Address
4     unsigned long vm_end;     // Logical End Address
5     struct vm_rg_struct *vm_freerg_list; // List of free holes for reuse
6     struct vm_area_struct *vm_next;      // Pointer to next VMA
7 };
```

Listing 6: Virtual Memory Area Structure

2. Secure User/Kernel Interface (`libstd.c`)

This wrapper demonstrates the protection mechanism. The user space only passes value arguments and the `pid`, never the raw PCB pointer, ensuring isolation.

```
1 int libsyscall(struct pcb_t *caller, uint32_t syscall_idx,
2               arg_t a1, arg_t a2, arg_t a3)
3 {
4     struct sc_regs regs;
5     regs.a1 = a1; regs.a2 = a2; regs.a3 = a3;
6
7     /* CRITICAL: Only caller->pid is passed to the kernel.
8        The Kernel must resolve the PCB internally. */
9     return syscall(caller->krnl, caller->pid, syscall_idx, &regs);
10 }
```

Listing 7: User Space System Call Wrapper

3. Memory Allocation Algorithm - `__alloc` (`libmem.c`)

This function executes in **Kernel Mode** after the PID has been validated. It implements the logic to search for available memory space within the VMA.

```
1 int __alloc(struct pcb_t *caller, int vmaid, int rgid, addr_t size, addr_t *
2             alloc_addr)
3 {
4     pthread_mutex_lock(&mmvm_lock);
5     struct vm_rg_struct rgnode;
```

```
6 // Check if we can reuse a free region
7 if (get_free_vmrg_area(caller, vmaid, size, &rgnode) == 0)
8 {
9     caller->mm->symrgtbl[rgid].rg_start = rgnode.rg_start;
10    caller->mm->symrgtbl[rgid].rg_end = rgnode.rg_end;
11    *alloc_addr = rgnode.rg_start;
12    pthread_mutex_unlock(&mmvm_lock);
13    return 0;
14 }
15
16 /* If get_free_vmrg_area FAILED, we must expand the heap */
17 struct vm_area_struct *cur_vma = get_vma_by_num(caller->mm, vmaid);
18 if (!cur_vma)
19 {
20     pthread_mutex_unlock(&mmvm_lock);
21     return -1;
22 }
23
24 addr_t inc_sz = PAGING_PAGE_ALIGNSZ(size);
25 // Usually aligning to page size is good practice, or use exact size if your
26 // inc_vma handles it.
27 // The provided skeleton had some ifdef logic, let's stick to standard alignment.
28
29 addr_t old_sbrk = cur_vma->sbrk;
30
31 /* SYSCALL to increase limit */
32 struct sc_regs regs;
33 regs.a1 = SYSMEM_INC_OP;
34 regs.a2 = vmaid;
35 regs.a3 = inc_sz; // Request expansion
36
37 // Note: In real syscall, we don't pass PCB, but here we simulate it via wrapper
38 // The wrapper 'syscall' takes (krnl, pid, nr, regs)
39 if (syscall(caller->krnl, caller->pid, 17, &regs) < 0)
40 {
41     pthread_mutex_unlock(&mmvm_lock);
42     return -1; // Failed to expand
43 }
44
45 /* Successful increase limit */
46 caller->mm->symrgtbl[rgid].rg_start = old_sbrk;
47 caller->mm->symrgtbl[rgid].rg_end = old_sbrk + size;
48 *alloc_addr = old_sbrk;
```



```
49 // The remaining space (inc_sz - size) should technically be added to free list
50 // to avoid internal fragmentation, but for this simple assignment we might skip
    it
51 // or add logic:
52 if (inc_sz > size)
53 {
54     struct vm_rg_struct *fragment = malloc(sizeof(struct vm_rg_struct));
55     fragment->rg_start = old_sbrk + size;
56     fragment->rg_end = old_sbrk + inc_sz;
57     fragment->rg_next = NULL;
58     enlist_vm_freerg_list(caller->mm, fragment);
59 }
60
61 pthread_mutex_unlock(&mmvm_lock);
62 return 0;
63 }
```

Listing 8: Kernel Space Allocation Routine

4. Memory Deallocation Algorithm - enlist_vm_freerg_list (libmem.c)

This function inserts a recently freed memory region back into the reuse list (free list), marking it available for future allocations.

```
1 int enlist_vm_freerg_list(struct mm_struct *mm, struct vm_rg_struct *rg_elmt)
2 {
3     struct vm_rg_struct *rg_node = mm->mmap->vm_freerg_list;
4
5     if (rg_elmt->rg_start >= rg_elmt->rg_end)
6         return -1;
7
8     if (rg_node != NULL)
9         rg_elmt->rg_next = rg_node;
10
11     /* Enlist the new region */
12     mm->mmap->vm_freerg_list = rg_elmt;
13
14     return 0;
15 }
```

Listing 9: Deallocation Logic

c. Algorithm Description Step 1: User-Kernel Transition (Protection Mechanism) To ensure memory safety (referencing *Silberschatz Chapter 2*), the allocation process begins with a System Call:

- **User Request:** The user program calls a library function (e.g., `malloc`), which triggers the wrapper `libsyscall`.

- **Trap to OS:** `libsycall` passes the Process ID (PID) and an operation code (such as `SYSMEM_INC_OP` in our allocation path, or `SYSMEM_MAP_OP` in other mapping scenarios) to the Kernel via a software interrupt/trap.
- **Authentication:** Upon entering Kernel Mode, the kernel uses the provided PID to traverse the process table and retrieve the correct `pcb_t` pointer. This prevents any user process from manipulating the memory of another process.

Step 2: VMA Management (`get_vma_by_num`)

Once in Kernel Mode, the system must identify which memory region the process is requesting.

- **Activity:** The function `get_vma_by_num` traverses the linked list `mm->mmap` to find the VMA whose `vm_id` equals the requested `vmid`.
- **Purpose:** In our implementation, there is exactly one user VMA per process (`vmid = 0`), which acts as the container for all user regions. The logical separation between different user “segments” (variables/objects) is handled by `vm_rg_struct` entries stored in `mm->symrgtbl[]`, rather than by multiple VMAs.

Step 3: Memory Allocation (`__alloc`)

The allocation logic prioritizes filling "holes" in the memory space before expanding it, ensuring resource efficiency.

1. **Reuse Strategy:** The system inspects the `vm_freerg_list` of the current VMA. This list contains memory regions that were previously freed. The algorithm applies a **First-Fit** strategy: it traverses the list, and if a node satisfying $rg_end - rg_start \geq size$ is found, the system either splits the node or claims the entire node for the new allocation.
2. **Expansion Strategy:** If the free list is empty or contains no suitable regions, the system performs a linear expansion. The new allocation address starts at the current `vm_end`, and the boundary is updated: $vm_end = vm_end + size$. This behavior is analogous to the Unix `sbrk()` system call.
3. **Registration:** Finally, the system updates the `symrgtbl` (Symbol Region Table) with the `rgid`. This maps a User-managed Region ID to the actual Kernel-managed logical address.

Step 4: Memory Deallocation (`enlist_vm_freerg_list`)

When a process frees memory, the region is not immediately physically erased but transitioned to an "Available" state.

- **Logic:** The function receives a memory region structure (`rg_elmt`).
- **Operation:** It inserts this node into the head of the singly linked list `vm_freerg_list`.
- **Significance:** This turns used memory into a new "hole". The next allocation request (via `__alloc`) will scan and potentially reuse this region, thereby minimizing external fragmentation of the logical address space.

Conclusion:

The Memory Management module at the Logical Level **implements** a memory management mechanism **predicated on** strict User/Kernel separation and an allocation strategy **utilizing** Free List reuse. This architecture

functions as the initial abstraction layer, establishing the logical address space before these addresses are translated into physical frames by the 5-Level Paging mechanism detailed in the subsequent section.

2.2.2 Multi-level paging

a. Implementation Details The following core components constitute the paging mechanism.

1. Bit Masking and Address Decoding (mm64.h) The system uses macros to "peel off" each layer of the virtual address.

```
1  /* PGD Index: Extract 9 bits from position 48 to 56 */
2  #define PAGING64_ADDR_PGD_MASK    GENMASK64(56, 48)
3  #define PAGING64_ADDR_PGD_LOBIT   48
4  #define PAGING64_ADDR_PGD(addr)   ((addr & PAGING64_ADDR_PGD_MASK) >>
   PAGING64_ADDR_PGD_LOBIT)
5
6  /* P4D -> PUD -> PMD: Shift the bit window downwards */
7  #define PAGING64_ADDR_P4D(addr)   ((addr & GENMASK64(47, 39)) >> 39)
8  #define PAGING64_ADDR_PUD(addr)   ((addr & GENMASK64(38, 30)) >> 30)
9  #define PAGING64_ADDR_PMD(addr)   ((addr & GENMASK64(29, 21)) >> 21)
10
11 /* PT Index (Leaf Level): Bits 20-12 */
12 #define PAGING64_ADDR_PT_MASK      GENMASK64(20, 12)
13 #define PAGING64_ADDR_PT(addr)     ((addr & PAGING64_ADDR_PT_MASK) >> 12)
14
15 /* Offset: Last 12 bits */
16 #define PAGING64_ADDR_OFFST(addr)  (addr & GENMASK64(11, 0))
```

Listing 10: Bit Manipulation Macros in mm64.h

Page Table Entry - PTE (mm64.h & mm64.c) The PTE stores the page state. Manipulating the PTE requires bitwise operations (`|`, `&`, `~`).

```
1  /* mm64.h - PTE Bit Definitions */
2  #define PAGING_PTE_PRESENT_MASK   BIT_ULL(0)   // Bit 0: Page is in RAM
3  #define PAGING_PTE_DIRTY_MASK     BIT_ULL(6)   // Bit 6: Page was written to
4  #define PAGING_PTE_SWAPPED_MASK   BIT_ULL(9)   // Bit 9: Page is swapped out
5  #define PAGING_PTE_FPN_MASK       GENMASK64(51, 12) // Bits 12-51: Frame Number
6
7  /* mm64.c - PTE Initialization Function */
8  int init_pte(addr_t *pte, int pre, addr_t fpn, int drt, int swp, int swptyp, addr_t
   swpoff) {
9      if (pre != 0) { // If page exists (Present)
10         if (swp == 0) { // Case 1: Page is in RAM
11             SETBIT(*pte, PAGING_PTE_PRESENT_MASK);
12             CLRBIT(*pte, PAGING_PTE_SWAPPED_MASK);
13             /* Write FPN to bits 12-51 */
14             SETVAL(*pte, fpn, PAGING_PTE_FPN_MASK, PAGING_PTE_FPN_LOBIT);
```

```
15     } else { // Case 2: Page is Swapped
16         SETBIT(*pte, PAGING_PTE_SWAPPED_MASK);
17         CLRBIT(*pte, PAGING_PTE_PRESENT_MASK);
18         /* Write Swap Offset */
19         SETVAL(*pte, swpoff, PAGING_PTE_SWPOFF_MASK, PAGING_PTE_SWPOFF_LOBIT);
20     }
21 }
22 return 0;
23 }
```

Listing 11: PTE Definition and Initialization

3. Page Walk Mechanism (mm64.c) The `get_page_table_entry` function is the core of translation, traversing 5 table levels.

```
1 uint64_t *get_page_table_entry(struct mm_struct *mm, addr_t addr, int alloc)
2 {
3     // Get Indices
4     int pgd_idx = PAGING64_ADDR_PGD(addr);
5     int p4d_idx = PAGING64_ADDR_P4D(addr);
6     int pud_idx = PAGING64_ADDR_PUD(addr);
7     int pmd_idx = PAGING64_ADDR_PMD(addr);
8     int pt_idx = PAGING64_ADDR_PT(addr);
9
10    // Check Root (PGD)
11    if (mm->pgd == NULL)
12    {
13        if (!alloc)
14            return NULL;
15        mm->pgd = alloc_aligned_table();
16
17        if (!mm->pgd)
18            return NULL;
19        memset(mm->pgd, 0, sizeof(struct pgd_t));
20    }
21
22    // Walk PGD -> P4D
23    struct p4d_t *p4d_table = get_next_level((uint64_t *)&mm->pgd->entries[pgd_idx],
24        alloc);
25    if (!p4d_table)
26        return NULL;
27
28    // Walk P4D -> PUD
29    struct pud_t *pud_table = get_next_level((uint64_t *)&p4d_table->entries[p4d_idx],
30        alloc);
31    if (!pud_table)
```

```
30     return NULL;
31
32     // Walk PUD -> PMD
33     struct pmd_t *pmd_table = get_next_level((uint64_t *)&pud_table->entries[pud_idx], alloc);
34     if (!pmd_table)
35         return NULL;
36
37     // Walk PMD -> PT
38     struct pt_t *pt_table = get_next_level((uint64_t *)&pmd_table->entries[pmd_idx], alloc);
39     if (!pt_table)
40         return NULL;
41
42     // Return pointer to the specific Page Table Entry (Leaf)
43     return (uint64_t *)&pt_table->entries[pt_idx];
44 }
```

Listing 12: Page Walk Logic (Simplified)

4. Hardware Simulation (mm-memphy.c) RAM is simulated not as a physical chip but as a BYTE array.

```
1  /* mm-memphy.c */
2  int MEMPHY_read(struct memphy_struct *mp, addr_t addr, BYTE *value) {
3      if (mp == NULL) return -1;
4      /* Direct access to storage array at index addr */
5      *value = mp->storage[addr];
6      return 0;
7  }
8
9  int MEMPHY_write(struct memphy_struct *mp, addr_t addr, BYTE data) {
10     if (mp == NULL) return -1;
11     mp->storage[addr] = data;
12     return 0;
13 }
```

Listing 13: Physical Memory Access

b. Algorithm Description The system executes address translation via a strict procedure, simulating the hardware MMU (Memory Management Unit).

Step 1: Address Decoding The computer views the virtual address not as a large integer, but as a collection of indices.

- **Input:** 64-bit Virtual Address (VA).
- **Algorithm:** Use Bitwise AND (&) with a Mask to isolate specific bits, then Bitwise Right Shift (>>) to normalize the value into an integer index.

- **Example:** To get the PGD Index, the system extracts the highest 9 bits (56-48). This value determines the position in the PGD table.

Step 2: Page Walk (Tree Traversal) This is the core algorithm. Conceptually, the page table is a 5-level tree.

1. **Root:** Start from `mm->pgd` (stored in the CR3 register on real CPUs).
2. **Traversal:**
 - Use PGD Index to select an entry in PGD. This entry contains the physical address of the P4D table.
 - Use P4D Index to select an entry in P4D. This points to the PUD table.
 - Continue similarly through PUD and PMD.
3. **Leaf:** At the final level (PT), use the PT Index to retrieve the **PTE**.

Logic Flow: `CR3 → [PGD] → [P4D] → [PUD] → [PMD] → [PT] → Frame`.

Step 3: PTE Handling & State Check Once the PTE is retrieved, the system checks the status bits:

- **Case A: Bit Present = 1 (Hit):** The page is in RAM. The system extracts the 40-bit FPN (Frame Page Number) from the PTE.
- **Case B: Bit Present = 0 & Swapped = 1 (Page Fault):** The data is not in RAM but in Swap. The system triggers `__mm_swap_page`: Find a free frame → Read data from Swap → Update PTE (Present=1, Swapped=0).
- **Case C: Bit Present = 0 & Swapped = 0 (Invalid):** Accessing unallocated memory → Segmentation Fault.

Step 4: Physical Address Composition After obtaining the frame number (FPN), the final physical address is calculated:

$$PhysicalAddress = (FPN \times PageSize) + Offset$$

In code, this is optimized using bitwise operations:

- `FPN << 12`: Left shift 12 bits (equivalent to multiplying by 4096).
- `| Offset`: Bitwise OR with Offset (lower 12 bits of VA) to combine them.

Step 5: Physical Access Finally, the Physical Address (PA) is passed to `mm-memphy.c`.

- The module treats RAM as a massive array `storage[]`.
- It accesses `storage[PA]` to read or write the actual data byte.

Conclusion The successful implementation of the 5-Level Paging mechanism demonstrates memory management at the finest granularity. The Page Walk algorithm, combined with precise Bitwise operations, ensures the system can translate addresses within the vast 64-bit space while transparently handling complex scenarios like Page Faults and Swapping.

3 Interprets the results of running tests

3.1 Scheduling

3.1.1 The result input file of scheduler (sched)

The following output describes the behavior of the scheduler over multiple time slots, showing how processes are loaded, dispatched, and processed by the CPU.

```
1 Time slot    0
2 ld_routine
3     Loaded a process at input/proc/p2s, PID: 1 PRI0: 0
4 Time slot    1
5     CPU 1: Dispatched process 1
6     Loaded a process at input/proc/p1s, PID: 2 PRI0: 1
7 Time slot    2
8     CPU 0: Dispatched process 2
9 Time slot    3
10    Loaded a process at input/proc/p2s, PID: 3 PRI0: 1
11 Time slot    4
12    Loaded a process at input/proc/p3s, PID: 4 PRI0: 0
13 Time slot    5
14    CPU 1: Put process 1 to run queue
15    CPU 1: Dispatched process 4
16 Time slot    6
17    CPU 0: Put process 2 to run queue
18    CPU 0: Dispatched process 1
19 Time slot    7
20 Time slot    8
21 Time slot    9
22    CPU 1: Put process 4 to run queue
23    CPU 1: Dispatched process 4
24 Time slot   10
25    CPU 0: Put process 1 to run queue
26    CPU 0: Dispatched process 1
27 Time slot   11
28 Time slot   12
29 Time slot   13
30    CPU 1: Put process 4 to run queue
31    CPU 1: Dispatched process 4
32 Time slot   14
33    CPU 0: Processed 1 has finished
34    CPU 0: Dispatched process 3
35 Time slot   15
36 Time slot   16
37    CPU 1: Processed 4 has finished
```

```
38         CPU 1: Dispatched process 2
39 Time slot 17
40 Time slot 18
41         CPU 0: Put process 3 to run queue
42         CPU 0: Dispatched process 3
43 Time slot 19
44 Time slot 20
45         CPU 1: Put process 2 to run queue
46         CPU 1: Dispatched process 2
47 Time slot 21
48 Time slot 22
49         CPU 0: Put process 3 to run queue
50         CPU 0: Dispatched process 3
51         CPU 1: Processed 2 has finished
52         CPU 1 stopped
53 Time slot 23
54 Time slot 24
55 Time slot 25
56 Time slot 26
57         CPU 0: Processed 3 has finished
58         CPU 0 stopped
```

Explanation of Output

The output demonstrates the behavior of the scheduler across multiple time slots and the corresponding actions performed by the CPU. Below is an analysis of key steps:

- **Time Slot 0 - 4 (Loading and Initial Dispatching):**

- *Loading:* Processes are loaded sequentially. Process 1 (Prio 0) loads at slot 0; Process 2 (Prio 1) at slot 1; Process 3 (Prio 1) at slot 3; and Process 4 (Prio 0) at slot 4.
- *Dispatching:* CPU 1 picks up Process 1 at slot 1. CPU 0 picks up Process 2 at slot 2. This shows the workload being distributed across both CPUs immediately as processes arrive.

- **Time Slot 5 - 6 (Preemption and Context Switching):**

- *Preemption:* At slot 5, Process 1 (on CPU 1) exhausts its time slice and is put back into the run queue. CPU 1 immediately dispatches the newly loaded Process 4 (Prio 0).
- *Process Migration:* At slot 6, Process 2 (on CPU 0) is preempted. CPU 0 then picks up Process 1 (which was previously on CPU 1). This clearly demonstrates *process migration* in the SMP environment.

- **Time Slot 9 - 13 (Re-dispatching):**

- *Round-Robin Behavior:* Processes 4 and 1 continue to execute. At slot 9 and 13, CPU 1 puts Process 4 back to the queue but immediately re-dispatches it (likely because it has the highest priority/quota at that moment). Similarly, CPU 0 re-dispatches Process 1 at slot 10.

• **Time Slot 14 - 16 (Process Completion):**

- *Completion:* Process 1 finishes execution on CPU 0 at slot 14. CPU 0 then switches to Process 3 (Prio 1).
- *Completion:* Process 4 finishes execution on CPU 1 at slot 16. CPU 1 then switches to Process 2 (Prio 1), which had been waiting since slot 6.

• **Time Slot 22 - 26 (Finalization):**

- *Termination:* Process 2 finishes at slot 22, and CPU 1 stops as there are no more processes for it. Process 3 continues on CPU 0 until slot 26, where it finishes, and CPU 0 finally stops.

3.1.2 Scheduling: Gantt chart

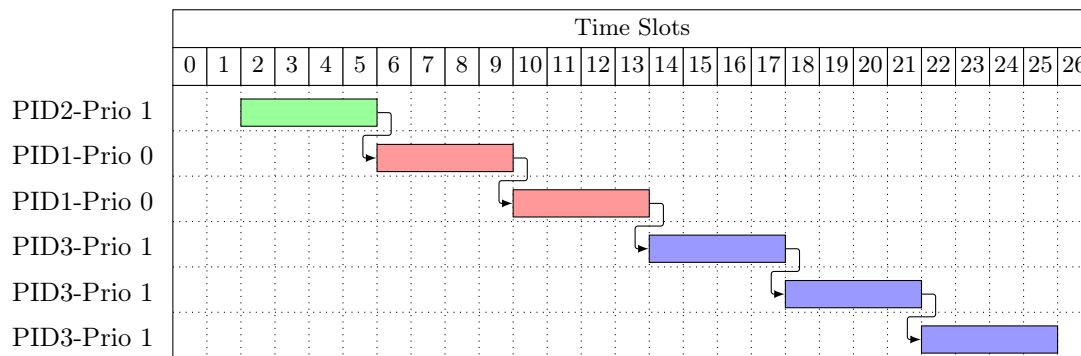


Figure 2: Gantt Chart for CPU 0

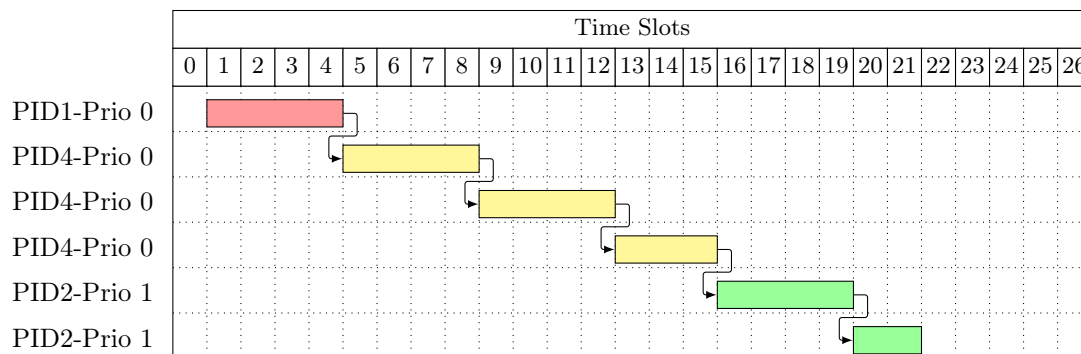


Figure 3: Gantt Chart for CPU 1

3.2 Memory

3.2.1 Result of the input file of Memory (os_0_mlq_paging)

```

1   Time slot    0
2 ld_routine

```



```
3      Loaded a process at input/proc/p0s, PID: 1 PRI0: 0
4 Time slot  1
5      CPU 1: Dispatched process  1
6      Loaded a process at input/proc/pls, PID: 2 PRI0: 15
7 Time slot  2
8      CPU 0: Dispatched process  2
9 liballoc:152
10 --- Page Table Dump ---
11 print_pgtbl:
12 PGD=0x739974001000 P4D=0x739974003000 PUD=0x739974005000 PMD=0x739974007000 PT=0
    x739974009000
13 0000000000000000: 0000000000000001
14 Time slot  3
15      Loaded a process at input/proc/pls, PID: 3 PRI0: 0
16 liballoc:152
17 --- Page Table Dump ---
18 print_pgtbl:
19 PGD=0x739974001000 P4D=0x739974003000 PUD=0x739974005000 PMD=0x739974007000 PT=0
    x739974009000
20 0000000000000000: 0000000000000001
21 Time slot  4
22 libfree:168
23 --- Page Table Dump ---
24 print_pgtbl:
25 PGD=0x739974001000 P4D=0x739974003000 PUD=0x739974005000 PMD=0x739974007000 PT=0
    x739974009000
26 0000000000000000: 0000000000000001
27 Time slot  5
28 liballoc:152
29 --- Page Table Dump ---
30 print_pgtbl:
31 PGD=0x739974001000 P4D=0x739974003000 PUD=0x739974005000 PMD=0x739974007000 PT=0
    x739974009000
32 0000000000000000: 0000000000000001
33 Time slot  6
34      Loaded a process at input/proc/p0s, PID: 4 PRI0: 0
35 libwrite:343
36 --- Page Table Dump ---
37 print_pgtbl:
38 PGD=0x739974001000 P4D=0x739974003000 PUD=0x739974005000 PMD=0x739974007000 PT=0
    x739974009000
39 0000000000000000: 0000000000000001
40 Time slot  7
41      CPU 1: Put process  1 to run queue
```



```
42      CPU 1: Dispatched process  3
43 Time slot  8
44      CPU 0: Put process  2 to run queue
45      CPU 0: Dispatched process  4
46 Time slot  9
47 liballoc:152
48 --- Page Table Dump ---
49 print_pgtbl:
50 PGD=0x73996c001000 P4D=0x73996c003000 PUD=0x73996c005000 PMD=0x73996c007000 PT=0
    x73996c009000
51 0000000000000000: 0000000000001001
52 Time slot 10
53 liballoc:152
54 --- Page Table Dump ---
55 print_pgtbl:
56 PGD=0x73996c001000 P4D=0x73996c003000 PUD=0x73996c005000 PMD=0x73996c007000 PT=0
    x73996c009000
57 0000000000000000: 0000000000001001
58 Time slot 11
59 libfree:168
60 --- Page Table Dump ---
61 print_pgtbl:
62 PGD=0x73996c001000 P4D=0x73996c003000 PUD=0x73996c005000 PMD=0x73996c007000 PT=0
    x73996c009000
63 0000000000000000: 0000000000001001
64 Time slot 12
65 liballoc:152
66 --- Page Table Dump ---
67 print_pgtbl:
68 PGD=0x73996c001000 P4D=0x73996c003000 PUD=0x73996c005000 PMD=0x73996c007000 PT=0
    x73996c009000
69 0000000000000000: 0000000000001001
70 Time slot 13
71 libwrite:343
72 --- Page Table Dump ---
73 print_pgtbl:
74 PGD=0x73996c001000 P4D=0x73996c003000 PUD=0x73996c005000 PMD=0x73996c007000 PT=0
    x73996c009000
75 0000000000000000: 0000000000001001
76      CPU 1: Put process  3 to run queue
77      CPU 1: Dispatched process  1
78 libread:304
79 Time slot 14
80      CPU 0: Put process  4 to run queue
```

```
81      CPU 0: Dispatched process 3
82 Time slot 15
83 Time slot 16
84 Time slot 17
85 Time slot 18
86      CPU 0: Processed 3 has finished
87      CPU 0: Dispatched process 4
88 libread:304
89 Time slot 19
90      CPU 1: Put process 1 to run queue
91      CPU 1: Dispatched process 1
92 libfree:168
93 --- Page Table Dump ---
94 print_pgtbl:
95 PGD=0x739974001000 P4D=0x739974003000 PUD=0x739974005000 PMD=0x739974007000 PT=0
    x739974009000
96 0000000000000000: 0000000000000001
97 Time slot 20
98 Time slot 21
99      CPU 1: Processed 1 has finished
100     CPU 1: Dispatched process 2
101 Time slot 22
102 Time slot 23
103 Time slot 24
104     CPU 0: Put process 4 to run queue
105     CPU 0: Dispatched process 4
106 libfree:168
107 --- Page Table Dump ---
108 print_pgtbl:
109 PGD=0x73996c001000 P4D=0x73996c003000 PUD=0x73996c005000 PMD=0x73996c007000 PT=0
    x73996c009000
110 0000000000000000: 00000000000001001
111 Time slot 25
112     CPU 1: Processed 2 has finished
113     CPU 1 stopped
114 Time slot 26
115     CPU 0: Processed 4 has finished
116     CPU 0 stopped
```

3.2.2 Explain the status of the memory allocation in heap and data segments

- **Heap Segment:** It is the memory area for initialization (Dynamic Allocation) while the program is running. It has no fixed size and can be expanded.

- **Data Segment:** In this simplified simulation model, the "Data Segment" refers to the actual byte values written into the allocated Heap memory frames. While the Heap defines the "container" (virtual addresses and size), the Data represents the "content" stored in the corresponding Physical Frames (RAM).

3.2.3 Chronological Analysis of Memory Allocation

The following analysis tracks the state of memory allocation across specific time slots where memory intervention occurs. Two distinct processes are identified by their Page Global Directory (PGD) addresses:

- **Process A (p0s):** PGD Address 0x739974001000
- **Process B (p1s):** PGD Address 0x73996c001000

Start Explain the status of the memory allocation in heap and data segments.

- **Time Slot 2: Initial Allocation (Process A)**
 - **Event:** liballoc:152 triggered by Process A.
 - **Action:** The process requests memory allocation. The OS expands the Heap and creates a new virtual memory region starting at Virtual Address 0.
 - **Page Table Dump:** 00...00: 00...01
 - **Analysis:** The Page Table Entry (PTE) value 0x1 indicates:
 - * **Present Bit (Bit 0):** 1 (Page is in RAM).
 - * **Frame Page Number (FPN):** 0 (Derived from the upper bits).

Process A is assigned **Physical Frame 0**.

- **Time Slot 4: Logical Deallocation (Process A)**
 - **Event:** libfree:168.
 - **Action:** Process A calls `free`.
 - **Heap Status:** The memory region is logically marked as "free" and added to the `vm_freerg_list` for future reuse.
- **Time Slot 5: Re-allocation (Process A)**
 - **Event:** liballoc:152.
 - **Action:** Process A requests allocation again.
 - **Heap Status:** The OS checks the `vm_freerg_list`, finds the region freed in Time Slot 4, and reuses it.
 - **Mapping:** The mapping remains 00...01 (Physical Frame 0).
- **Time Slot 6: Data Write (Process A)**
 - **Event:** libwrite:343.
 - **Action:** Data is written to the allocated address.

- **Data Status:** The actual value is stored in **Physical Frame 0**.
- **Time Slot 9: Allocation for New Process (Process B)**
 - **Event:** liballoc:152.
 - **Context Switch:** The PGD changes to 0x73996c001000, indicating a context switch to Process B.
 - **Page Table Dump:** 00...00: 00...1001
 - **Analysis:** The PTE value 0x1001 indicates:
 - * **Present Bit:** 1.
 - * **FPN:** 1 (The value 0x1000 at bit offset 12 corresponds to index 1).

Since Frame 0 is occupied by Process A, Process B is assigned the next available resource, **Physical Frame 1**.
- **Time Slot 11 & 13: Free and Write (Process B)**
 - **Time Slot 11 (Free):** Process B frees its memory region. The region enters Process B's private free list.
 - **Time Slot 13 (Write):** Process B writes data. This data is stored in **Physical Frame 1**.
- **Time Slot 19 & 24: Final Deallocation** Both processes perform their final libfree operations (Slot 19 for Process A, Slot 24 for Process B), releasing their respective logical memory regions.

3.2.4 Summary of Result Flow

The table below summarizes the correlation between Time Slots, Processes, and Physical Memory status.

| Slot | PGD (Process) | Event | Physical Frame | Logical Heap Status |
|------|------------------|-------|----------------|------------------------------|
| 2 | ...7400 (Proc A) | Alloc | Frame 0 | New Region created |
| 4 | ...7400 (Proc A) | Free | Frame 0 | Region moved to Free List |
| 5 | ...7400 (Proc A) | Alloc | Frame 0 | Region reused from Free List |
| 6 | ...7400 (Proc A) | Write | Frame 0 | Data written to Frame 0 |
| 9 | ...6c00 (Proc B) | Alloc | Frame 1 | New Region created |
| 11 | ...6c00 (Proc B) | Free | Frame 1 | Region moved to Free List |
| 13 | ...6c00 (Proc B) | Write | Frame 1 | Data written to Frame 1 |

Table 1: Memory Allocation Flow Summary

4 Answer question

4.1 Question 1

Question: What is the mechanism to pass a complex argument to a system call using the limited registers?

Answer:

The mechanism used to pass a complex argument to a system call is Indirection via pointers. Since CPU registers are very limited in number and size, they cannot hold large or complex data directly. Instead, the user-space program places the data in its own memory and passes only the memory address (pointer) of that data in a register. When the system call is invoked, the kernel receive the pointer from the register. But this method have a problem that the passed address maybe wrong or bad address (pointing to kernel memory or unmapped memory) leading to the kernel could crash. To handle this problem, the kernel check if the address provided is actually within the User space memory area. Then the kernel copy the data from User space into Kernel Space buffer, the Kernel perform on its own safe copy of the data.

4.2 Question 2

Question: Which mechanisms does the operating system use to manage system calls that become unresponsive?

Answer:

Operating systems use a multi-layered approach to manage system calls that become unresponsive. The mechanisms used to manage this fall into three categories:

- **Interruption Mechanisms:** when the system call blocks, the operating system provides ways for the user or the application to forcefully break. In Linux most of blocking system calls place the process in a task interruptible state, when the process receives a signal like Ctrl + C, the kernel prematurely wakes the process.
- **Timeout/Watchdogs:** The kernel monitors itself to ensure that a single stuck system call does not freeze the entire computer. Many blocking system calls accept a timeout parameter. The kernel sets an internal timer; if the condition is not satisfied before the timer expires, the call returns with an error. In addition, the operating system runs a background kernel thread that scans all process in every constant time (120 seconds in Linux), if it find a process stuck in Uninterruptible Sleep (D state) for more than this duration, it logs a "Hung Task" error to the kernel log with a stack trace.
- **Asynchronous Avoidance:** Modern operating systems encourage developers to avoid "blocking" system calls entirely to prevent unresponsiveness. Applications can set file descriptors to non-blocking mode (O_NONBLOCK). Its mean that if a system call would freeze, it returns immediately with an error instead of waiting. In addition, there are some API (Application Programming Interface) like io_uring, epoll,.. in Linux that allow applications to submit system calls to a queue and retrieve results later, ensuring the main application thread never actually block.

4.3 Question 3

Question: Evaluate the comparative advantages of the scheduling algorithm implemented in this assignment in relation to other scheduling algorithms you have learned. What is the complexity of this algorithm?

Answer:

The scheduling algorithm implemented in this assignment is a **Non-Feedback Multilevel Queue (MLQ)** augmented with a **Slot-based Round Robin** mechanism for inter-queue dispatching. The design employs a static array of priority queues, where each priority level is assigned a specific execution quota (slot). Based on the theoretical framework provided in *Operating System Concepts (10th Edition)* and the specific logic in `sched.c`, we present a detailed comparative evaluation and complexity analysis below.

1. Comparative Evaluation against Standard Algorithms

a. Comparison with First-Come, First-Served (FCFS)

- **Theoretical Background:** FCFS is the most basic non-preemptive scheduling algorithm. The CPU is allocated to the process at the head of the ready queue. The critical drawback identified in operating system theory is the "*Convoy Effect*". This phenomenon occurs when a CPU-intensive process occupies the processor for an extended period, causing all subsequent I/O-bound processes to wait in the ready queue. This leads to extremely low I/O device utilization and a high Average Waiting Time.
- **Assignment Implementation Analysis:** The implemented algorithm in the assignment is inherently **Preemptive**.
 - *Mechanism:* The simulation framework in `os.c` enforces a `time_slot`. Even if a process has not completed its CPU burst, it is forced to yield the CPU and is placed back into the `mlq_ready_queue` via `put_mlq_proc()`.
 - *Code Evidence:* In `sched.c`, the scheduler manages `slot[i]`. When a process consumes its time slice, the scheduler is invoked again to pick a new candidate.
- **Comparative Advantage:** Our Slot-based MLQ algorithm completely eliminates the Convoy Effect. By forcing time-sharing (Round Robin within queues), short processes and I/O-bound processes are guaranteed frequent CPU access. This significantly improves the system's **Responsiveness**, making it suitable for interactive environments, whereas FCFS is strictly limited to Batch Systems.

b. Comparison with Shortest-Job-First (SJF)

- **Theoretical Background (Optimality vs. Prediction):** SJF (or more accurately, *Shortest-Next-CPU-Burst*) is theoretically proven to yield the minimum average waiting time. However, it suffers from a fundamental implementation flaw: the OS cannot know the length of the next CPU burst (τ_{n+1}) in advance.
To implement SJF, the system must rely on **Exponential Averaging** to predict the future based on past behavior, using the formula:

$$\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n \quad (2)$$

Where t_n is the most recent burst and τ_n is the past history. The parameter α ($0 \leq \alpha \leq 1$) controls the relative weight of recent versus past history:

- **Properties of α (Weighting):** The expansion of the formula ($\tau_{n+1} = \alpha t_n + (1 - \alpha)\tau_n + \dots$) shows that the weight of each past burst decays exponentially.
- **Scenario Analysis (When is it bad?):**
 - * If $\alpha \approx 0$ (History dominates): The system is sluggish. If a process changes behavior (e.g., from I/O-bound to CPU-bound), the scheduler fails to adapt, treating a long job as a short one.
 - * If $\alpha \approx 1$ (Recent dominates): The system overreacts to transient spikes. A single unusually long burst (an outlier) will cause the scheduler to penalize the process significantly in the next turn, even if the long burst was a one-time event.

Thus, while optimal in theory, SJF is computationally expensive (floating-point math in kernel) and prone to prediction errors when process behavior is unstable.

- **Assignment Implementation Analysis (Static Determinism):** Our implemented algorithm prioritizes **Determinism** and **Feasibility** over theoretical optimality by avoiding runtime predictions entirely.
 - *Mechanism:* Instead of guessing τ_{n+1} , the system relies on the explicit `proc->prio` value defined in the PCB. This value represents the "importance" or "class" of a process rather than its estimated burst length.
 - *Code Evidence:* The selection logic in `sched.c` (`get_mlq_proc`) iterates through `mlq_ready_queue[i]` using simple integer comparisons. This avoids the overhead of calculating exponential averages ($O(1)$ arithmetic vs. floating-point operations).
- **Comparative Evaluation:** Unlike SJF, which can degrade if α is poorly tuned or if processes exhibit erratic burst patterns (e.g., mixed CPU/IO workloads), our Slot-based MLQ guarantees that high-priority processes always run first, while the slot mechanism prevents the starvation that strict priority scheduling usually suffers from.

c. Comparison with Pure Priority Scheduling

- **Theoretical Background:** In standard Priority Scheduling, the CPU is always allocated to the highest priority process. The fatal flaw of this design is "*Indefinite Blocking*" or "*Starvation*". If a steady stream of high-priority processes arrives, low-priority processes may never execute. The standard solution is "Aging" (gradually increasing the priority of waiting processes), which adds complexity to the system.
- **Assignment Implementation Analysis (The Slot Mechanism):** This is the most innovative aspect of the assignment's design. It solves starvation **without** changing process priorities (Aging).

– *Code Logic:*

```
1 // sched.c : get_mlq_proc()
2 if (slot[i] > 0) {
3     proc = dequeue(...);
4     slot[i]--; // Consume budget
```

```
5     break;          // Dispatch immediate
6 } else {
7     slot[i] = MAX_PRIO - i; // Refill budget
8     continue; // SKIP this queue, force check lower priority
9 }
```

- *Operational Theory*: The system assigns a "CPU Budget" to each priority level based on the formula $Budget = MAX_PRIO - Prio$.
- *Example*: Priority 0 gets 140 slots. Priority 1 gets 139 slots. Even if Priority 0 has infinite tasks, after 140 dispatches, the `slot[0]` becomes 0. The scheduler is mathematically forced to skip Priority 0 and serve Priority 1.

- **Comparative Advantage**: The implementation achieves **Starvation Freedom** through a "Share-based" approach. It guarantees that every priority level, no matter how low, eventually receives a turn when higher levels exhaust their quotas. This is a robust and fairness-oriented improvement over standard Priority Scheduling.

d. Comparison with Standard Multilevel Queue (Partitioned)

- **Theoretical Background**: Standard MLQ partitions the ready queue into distinct groups, typically "Foreground" (Interactive) and "Background" (Batch). The scheduling between queues is often Fixed-Priority (leading to starvation) or Time-Slicing with coarse granularity (e.g., 80% CPU for Foreground, 20% for Background).
- **Assignment Implementation Analysis**: Our implementation offers a **Fine-Grained Proportional Share**.
 - *Granularity*: Instead of just 2 or 3 partitions, the system supports `MAX_PRIO` (140) distinct levels.
 - *Distribution*: The slot allocation decreases linearly (140, 139, ..., 1).
- **Comparative Advantage**: This linear distribution creates a smoother degradation of service. It allows the OS to differentiate processes with high precision (e.g., a process with priority 10 is slightly favored over 11), providing more flexibility than the rigid Foreground/Background binary of standard MLQ.

e. Comparison with Multilevel Feedback Queue (MLFQ)

- **Theoretical Background**: MLFQ is considered the most general and powerful scheduling algorithm. Its key feature is *Feedback*: processes move between queues. A process using too much CPU moves down; a process waiting too long moves up (Aging). This allows the OS to learn the nature of the process (I/O-bound vs CPU-bound) dynamically.
- **Assignment Implementation Limitation**: The assignment specification explicitly states a non-feedback design.
 - *Code Evidence*: In `put_mlq_proc()`, a process is always enqueued back to `mlq_ready_queue[proc->prio]`. The `prio` field is never modified by the scheduler.
- **Disadvantage/Trade-off**: The lack of feedback makes the system less adaptive. If a high-priority process becomes CPU-intensive, it continues to enjoy a large slot quota (140 slots),

potentially degrading overall system throughput compared to an MLFQ which would demote such a process. However, this simplifies the implementation significantly and reduces the overhead of priority recalculation.

2. Algorithmic Complexity Analysis

The performance of the scheduler is determined by the operations on the `queue_t` structure (array-based) and the search logic in `sched.c`. Let N be the number of processes in a specific queue, and K be the number of priority levels ($K = \text{MAX_PRIO} = 140$).

- **Enqueue Operation ($O(1)$):**

- *Code Reference:* `queue.c: enqueue()`.
- *Analysis:* The function inserts a process at the tail of the array: `q->proc[q->size] = proc`. This requires direct index access and incrementing the size counter.
- *Complexity:* Constant time, $O(1)$.

- **Dequeue Operation ($O(N)$):**

- *Code Reference:* `queue.c: dequeue()`.
- *Analysis:* The scheduler retrieves the process at the head (index 0). Because the queue is implemented as a contiguous memory block (array), removing the first element creates a "gap". The loop `for (i = 0; i < q->size - 1; i++)` is required to shift all remaining $N-1$ elements to the left to preserve data contiguity.
- *Complexity:* Linear time relative to queue size, $O(N)$.

- **Scheduler Selection Logic ($O(K + N)$):**

- *Code Reference:* `sched.c: get_mlq_proc()`.
- *Analysis:* The routine employs a linear search across priority levels: `for (int i = 0; i < MAX_PRIO; i++)`. In the worst-case scenario (e.g., when the system is idle or all high-priority slots are exhausted), the loop iterates K times. Upon finding a valid queue, it calls `dequeue()`.
- *Calculation:* $\text{Cost} = \text{Search_Cost} + \text{Retrieval_Cost} = O(K) + O(N)$.
- **Total Complexity:** $O(K + N)$.

4.4 Question 4

Question: Discuss the architectural and operational implications of implementing multiple memory segments in this simple OS.

Answer:

Implementing multiple memory segments fundamentally structures how the OS manages the virtual address space on top of the physical paging system. The implications can be divided into architectural changes and operational behaviors.

1. Architectural Implications

From an architectural perspective, this adds a logical organization layer above the raw paging mechanism:

- a. **Data Structure Transformation:** In the design, a segment is represented by a `vm_area_struct` inside the per-process `mm_struct`. Implementing multiple segments means that instead of a single VMA (current `vm_id = 0` for "data/heap"), `mm->mmap` becomes a linked list of VMAs. For example:
- `vm_id = 0`: Code / Read-only data
 - `vm_id = 1`: Heap
 - `vm_id = 2`: Stack
 - `vm_id = 3`: Shared regions, etc.
- b. **Logical Abstraction Layer:** Architecturally, this creates a logical layer on top of the 5-level page tables:
- The **5-level paging** (PGD–P4D–PUD–PMD–PT) remains unchanged and is solely responsible for translating virtual page numbers to physical frames.
 - The **Segmentation layer** decides *which* virtual ranges exist and *what* they are used for. Each VMA tracks its virtual interval [`vm_start`, `vm_end`), current growth point `sbrk`, and a free-region list for local holes.
- c. **Structured Extensibility:** This makes the address space structured and extensible. Different regions can later be assigned different protection policies (e.g., code as read-only, stack as non-executable, or shared segments mapped into multiple processes) while still sharing the same page table root (`mm->pgd`).

2. Operational Implications

On the operational side, supporting multiple segments changes the workflow of memory operations:

- a. **Segment-Aware Operations:** All core memory functions (`__alloc`, `__free`, `__read`, `__write`) become segment-aware. They take a `vmaid` parameter to select the correct VMA via `get_vma_by_num(mm, vmaid)` before performing any work. The allocator only searches and updates the free-region list of that specific segment, and `inc_vma_limit` grows only that targeted VMA when more virtual space is needed.
- b. **Bookkeeping vs. Isolation:** There is a trade-off involving increased bookkeeping for better isolation. Every time the kernel extends a segment, it must check for overlaps with other VMAs (`validate_overlap_vm_area`) before calling `vm_map_ram` to map new pages. While this adds overhead (walking the VMA list, extra validation), it guarantees that distinct segments (e.g., Heap vs. Stack) do not collide and grow independently.
- c. **Contextual Paging and Swapping:** While low-level page-table operations in `mm64.c` (allocating frames, filling PTEs, swapping victim pages) remain structurally the same, the context of each page becomes clearer. The OS knows whether a page belongs to code, heap, or stack based on the VMA it falls into. This facilitates segment-specific policies (e.g., preventing code pages from being swapped out or prioritizing stack pages) without altering the core paging algorithms.

In summary, while multiple segments slightly increase complexity in the memory manager, the OS gains a much more organized virtual address space and a natural architectural hook to implement protection and per-segment policies on top of the underlying 5-level paging mechanism.

4.5 Question 5

Question: What will happen if we divide the address to more than 2 levels in the paging memory management system?

Answer:

Based on the Simple OS implementation, dividing the address into more than 2 levels fundamentally changes the paging structure and introduces several important consequences:

1. Current 2-Level Implementation

In the Simple OS's 2-level paging system (defined in `common.h`), the 20-bit address space is divided as follows:

- **First level (Segment):** 5 bits (`FIRST_LV_LEN`) = 32 entries
- **Second level (Page):** 5 bits (`SECOND_LV_LEN`) = 32 entries
- **Offset:** 10 bits (`OFFSET_LEN`) = 1KB pages

This creates a hierarchical structure with `page_table_t` containing pointers to `trans_table_t` entries, requiring 2 memory accesses for address translation. The total virtual address space is limited to 1 MB (2^{20} bytes).

```
1 // common.h
2 /* Address geometry (two-level page tables):
3  * ADDRESS_SIZE: total VA bits.
4  * OFFSET_LEN:   bits within a page (page size = 1 << OFFSET_LEN).
5  * FIRST_LV_LEN: bits used by first-level index (L1).
6  * SECOND_LV_LEN: bits used by second-level index (L2).
7  * SEGMENT_LEN/PAGE_LEN: aliases to L1/L2 lengths.
8  */
9 #define ADDRESS_SIZE 20
10 #define OFFSET_LEN 10
11 #define FIRST_LV_LEN 5
12 #define SECOND_LV_LEN 5
13 #define SEGMENT_LEN FIRST_LV_LEN
14 #define PAGE_LEN SECOND_LV_LEN
```

Listing 14: 20-Bit Address Space

2. 5-Level Paging (64-bit Implementation)

When moving to 5-level paging in the 64-bit mode (`mm64.h` and `mm64.c`), the address is divided into 5 distinct levels:

- **Level 5 - PGD (Page Global Directory):** bits 56-48 (9 bits = 512 entries)
- **Level 4 - P4D (Page Level 4 Directory):** bits 47-39 (9 bits = 512 entries)
- **Level 3 - PUD (Page Upper Directory):** bits 38-30 (9 bits = 512 entries)

- **Level 2 - PMD (Page Middle Directory):** bits 29-21 (9 bits = 512 entries)
- **Level 1 - PT (Page Table):** bits 20-12 (9 bits = 512 entries)
- **Offset:** bits 11-0 (12 bits = 4KB pages)

This expands the virtual address space from 1 MB (20-bit) to 128 PiB (57-bit canonical address space), as implemented in `mm64.c`.

```
1 // mm64.h
2 #define MM64_BITS_PER_LONG 64
3
4 #define PAGING64_CPU_BUS_WIDTH 57 /* 57-bit canonical VA space */
5 #define PAGING64_PAGESZ 4096 /* 4KB or 12-bits PAGE NUMBER */
6
7 // ...
8
9 /* OFFSET */
10 #define PAGING64_ADDR_OFFST_HIBIT 11
11 #define PAGING64_ADDR_OFFST_LOBIT 0
12
13 /* PT */
14 #define PAGING64_ADDR_PT_HIBIT 20
15 #define PAGING64_ADDR_PT_LOBIT 12
16 #define PAGING64_ADDR_PT_SHIFT 12
17
18 /* PMD */
19 #define PAGING64_ADDR_PMD_HIBIT 29
20 #define PAGING64_ADDR_PMD_LOBIT 21
21
22 /* PUD */
23 #define PAGING64_ADDR_PUD_HIBIT 38
24 #define PAGING64_ADDR_PUD_LOBIT 30
25
26 /* P4D */
27 #define PAGING64_ADDR_P4D_HIBIT 47
28 #define PAGING64_ADDR_P4D_LOBIT 39
29
30 /* PGD */
31 #define PAGING64_ADDR_PGD_HIBIT 56
32 #define PAGING64_ADDR_PGD_LOBIT 48
```

Listing 15: 64-Bit 5 Level Paging

3. Key Consequences of Multi-Level Paging

- Address Space Expansion:** The 5-level scheme supports up to 128 petabytes of virtual memory, compared to the 2-level system's 1 MB limit. The function `get_pd_from_address()` in `mm64.c` (lines

149-157) extracts all 5 directory indices from a single address using bit masks and shifts, enabling this massive expansion. Each level uses 9 bits, allowing 512 entries per directory table.

```
1 // mm64.c
2 /*
3  * get_pd_from_address - Parse address to 5 page directory level
4  */
5 int get_pd_from_address(addr_t addr, addr_t *pgd, addr_t *p4d, addr_t *pud,
6   addr_t *pmd, addr_t *pt)
7 {
8     *pgd = (addr & PAGING64_ADDR_PGD_MASK) >> PAGING64_ADDR_PGD_LOBIT;
9     *p4d = (addr & PAGING64_ADDR_P4D_MASK) >> PAGING64_ADDR_P4D_LOBIT;
10    *pud = (addr & PAGING64_ADDR_PUD_MASK) >> PAGING64_ADDR_PUD_LOBIT;
11    *pmd = (addr & PAGING64_ADDR_PMD_MASK) >> PAGING64_ADDR_PMD_LOBIT;
12    *pt = (addr & PAGING64_ADDR_PT_MASK) >> PAGING64_ADDR_PT_LOBIT;
13    return 0;
14 }
```

Listing 16: get_pd_from_address() in mm64.c

b. Increased Translation Overhead: Each address translation requires traversing all 5 levels (PGD → P4D → PUD → PMD → PT), requiring 5 memory accesses instead of 2. The implementation in get_page_table_entry() (Listing 12) demonstrates this hierarchical traversal. Modern CPUs use Translation Lookaside Buffers (TLB) to cache recent translations, mitigating this overhead in practice.

c. Sparse Allocation and Memory Efficiency: The implementation uses **on-demand allocation** through the get_next_level() function (lines 72-96 of mm64.c). Page tables are allocated only when actually needed, not for the entire address space. This sparse allocation strategy is critical because:

- Theoretical maximum page table size would be $512^5 = 35$ trillion entries (unfeasible)
- Actual memory usage is much smaller since only used address ranges allocate tables
- Each table is 4KB (512 entries × 8 bytes), allocated via alloc_aligned_table()
- The mm_struct->pgd pointer starts as NULL and is allocated on first use

```
1 // mm64.c
2 /*
3  * Helper: Get or Allocate Next Level (Directory Navigation)
4  */
5 void *get_next_level(uint64_t *entry, int alloc)
6 {
7     /* Check Real Hardware Present Bit (Bit 0) */
8     if (*entry & X86_PRESENT)
9     {
10         /* Mask off the bottom 12 bits to get the pointer */
11     }
```

```
11         return (void *)(*entry & X86_ADDR_MASK);
12     }
13
14     if (!alloc)
15         return NULL;
16
17     /* Allocate 4KB Aligned Table */
18     void *new_table = alloc_aligned_table();
19     if (new_table == NULL)
20         return NULL;
21
22     /* * Write the Entry:
23      * Pointer (Already 4KB aligned, so bottom 12 bits are 0)
24      * OR with X86 Flags (Present + Write + User)
25      */
26     *entry = (uint64_t)new_table | X86_PRESENT | X86_WRITE | X86_USER;
27
28     return new_table;
29 }
```

Listing 17: get_next_level() in mm64.c

d. Page Table Entry Operations Complexity: All PTE operations must traverse the full 5-level hierarchy:

- pte_set_fpn() (lines 192-208): Walks 5 levels to set frame number

```
1 // mm64.c
2 /*
3  * pte_set_fpn - Set PTE entry for on-line page
4  */
5 int pte_set_fpn(struct pcb_t *caller, addr_t pgn, addr_t fpn)
6 {
7     addr_t addr = pgn << PAGING64_ADDR_PT_SHIFT;
8     uint64_t *pte = get_page_table_entry(caller->mm, addr, 1);
9     if (pte == NULL)
10         return -1;
11
12     *pte = 0; // Clear old data
13
14     // Mark page present in RAM
15     SETBIT(*pte, PAGING_PTE_PRESENT_MASK);
16     CLRBIT(*pte, PAGING_PTE_SWAPPED_MASK);
17
18     SETVAL(*pte, fpn, PAGING_PTE_FPN_MASK, PAGING_PTE_FPN_LOBIT);
19 }
```



```
20     return 0;
21 }
```

Listing 18: pte_set_fpn() in mm64.c

- pte_set_swap() (lines 171-187): Walks 5 levels to set swap information

```
1  // mm64.c
2  /*
3  * pte_set_swap - Set PTE entry for swapped page
4  */
5  int pte_set_swap(struct pcb_t *caller, addr_t pgn, int swptyp, addr_t
    swpoff)
6  {
7      addr_t addr = pgn << PAGING64_ADDR_PT_SHIFT;
8      uint64_t *pte = get_page_table_entry(caller->mm, addr, 1);
9      if (pte == NULL)
10         return -1;
11
12     *pte = 0; // Clear old data
13
14     // Swapped pages keep Present=0 and only carry swap metadata
15     SETBIT(*pte, PAGING_PTE_SWAPPED_MASK);
16
17     SETVAL(*pte, swptyp, PAGING_PTE_SWPTYP_MASK, PAGING_PTE_SWPTYP_LOBIT)
18         ;
19     SETVAL(*pte, swpoff, PAGING_PTE_SWPOFF_MASK, PAGING_PTE_SWPOFF_LOBIT)
20         ;
21
22     return 0;
23 }
```

Listing 19: pte_set_swap() in mm64.c

- pte_get_entry() (lines 211-218): Walks 5 levels to read PTE

```
1  // mm64.c
2  /* Get PTE page table entry */
3  uint64_t pte_get_entry(struct pcb_t *caller, addr_t pgn)
4  {
5      addr_t addr = pgn << PAGING64_ADDR_PT_SHIFT;
6      uint64_t *pte = get_page_table_entry(caller->mm, addr, 0);
7      if (pte == NULL)
8         return 0;
9      return *pte;
10 }
```

Listing 20: pte_get_entry() in mm64.c

- `vmap_pgd_memset()` (lines 234-246): Must walk 5 levels for each page mapped

```
1 // mm64.c
2 /*
3  * vmap_pgd_memset - map a range of page at aligned address
4  */
5 int vmap_pgd_memset(struct pcb_t *caller, addr_t addr, int pgnum)
6 {
7     for (int i = 0; i < pgnum; i++)
8     {
9         addr_t current_addr = addr + (i * PAGING64_PAGESZ);
10        uint64_t *pte = get_page_table_entry(caller->mm, current_addr, 1)
11        ;
12        if (pte)
13        {
14            *pte = 0; // Clear the LEAF entry (Not Present)
15        }
16    }
17    return 0;
18 }
```

Listing 21: `vmap_pgd_memset()` in `mm64.c`

Each operation calls `get_page_table_entry()` which performs the full traversal, allocating intermediate directories as needed.

- e. Implementation Complexity:** The 5-level implementation is significantly more complex than the 2-level system:

- `init_mm()` (lines 358-386 in `mm64.c`): Initializes the root PGD pointer as NULL, with VM area structures

```
1 // mm64.c
2 /* init_mm */
3 int init_mm(struct mm_struct *mm, struct pcb_t *caller)
4 {
5     struct vm_area_struct *vma0 = malloc(sizeof(struct vm_area_struct));
6     mm->pgd = NULL;
7
8     vma0->vm_id = 0;
9     vma0->vm_start = 0;
10    vma0->vm_end = vma0->vm_start;
11    vma0->sbrk = vma0->vm_start;
12    vma0->vm_freerg_list = NULL;
13
14    struct vm_rg_struct *first_rg = init_vm_rg(vma0->vm_start, vma0->
15        vm_end);
16    enlist_vm_rg_node(&vma0->vm_freerg_list, first_rg);
17 }
```

```
16
17     vma0->vm_next = NULL;
18     vma0->vm_mm = mm;
19     mm->mmap = vma0;
20
21     for (int i = 0; i < PAGING_MAX_SYMTBL_SZ; i++)
22     {
23         mm->symrgtbl[i].rg_start = 0;
24         mm->symrgtbl[i].rg_end = 0;
25         mm->symrgtbl[i].rg_next = NULL;
26     }
27
28     mm->fifo_pgn = NULL;
29
30     return 0;
31 }
```

Listing 22: init_mm() in mm64.c

- `get_page_table_entry()`: Implements recursive tree traversal with dynamic allocation (Listing 12)
- `print_pgtbl_recursive()` (lines 546-595 in mm64.c): Recursive function to traverse and print all 5 levels
- Address extraction in `get_pd_from_address()`: Complex bit manipulation across 5 levels (Listing 16)

The code structure shows separate directory types (`pgd_t`, `p4d_t`, `pud_t`, `pmd_t`, `pt_t`) defined in `os-mm.h`, each containing 512 entries.

```
1 // os-mm.h
2 #define MM_PAGING
3 #define PAGING_MAX_MMSWP 4 /* max number of supported swapped space */
4 #define PAGING_MAX_SYMTBL_SZ 30
5
6 /* 64-bit Paging Constants */
7 #define PAGING64_ENTRY_COUNT 512 // 9 bits per level = 512 entries
8
9 /*
10 * @bkssystem: in long address mode of 64bit or original 32bit
11 * the address type need to be redefined
12 */
13
14 #ifdef MM64
15 typedef uint64_t addr_t;
16 typedef char BYTE;
17
```

```

18 // Page Table (Level 1) - Points to Physical Frame
19 struct pt_t
20 {
21     uint64_t entries[PAGING64_ENTRY_COUNT];
22 };
23
24 // Page Middle Directory (Level 2) - Points to PT
25 struct pmd_t
26 {
27     uint64_t entries[PAGING64_ENTRY_COUNT];
28 };
29
30 // Page Upper Directory (Level 3) - Points to PMD
31 struct pud_t
32 {
33     uint64_t entries[PAGING64_ENTRY_COUNT];
34 };
35
36 // Page Level-4 Directory (Level 4) - Points to PUD
37 struct p4d_t
38 {
39     uint64_t entries[PAGING64_ENTRY_COUNT];
40 };
41
42 // Page Global Directory (Level 5) - Points to P4D
43 struct pgd_t
44 {
45     uint64_t entries[PAGING64_ENTRY_COUNT];
46 };

```

Listing 23: Directory Structs defined in os-mm.h

f. Storage Space Trade-offs: The implementation demonstrates careful design decisions to balance:

- **Traversal time:** 5 memory accesses per translation (mitigated by TLB)
- **Memory accesses:** Each level requires one access, but sparse allocation reduces total accesses
- **Storage space:** Only allocated tables consume memory, not theoretical maximum

The `vmap_page_range()` function (lines 251-278 of `mm64.c`) shows how pages are mapped efficiently, tracking mapped pages in the `fifo_pgn` list for replacement algorithms.

```

1 /*
2  * vmap_page_range - map a range of page at aligned address
3  */
4 addr_t vmap_page_range(struct pcb_t *caller,
5 addr_t addr,
6 int pgnum,

```

```
7 struct framephy_struct *frames,
8 struct vm_rg_struct *ret_rg)
9 {
10     struct framephy_struct *fpit = frames;
11     int pgit = 0;
12     addr_t pgn;
13
14     ret_rg->rg_start = addr;
15     ret_rg->rg_end = addr + (pgnum * PAGING64_PAGESZ);
16
17     for (pgit = 0; pgit < pgnum; pgit++)
18     {
19         if (fpit == NULL)
20             return -1;
21
22         pgn = (addr + (pgit * PAGING64_PAGESZ)) >> PAGING64_ADDR_PT_SHIFT;
23
24         pte_set_fpn(caller, pgn, fpit->fpn);
25         enlist_pgn_node(&caller->mm->fifo_pgn, pgn);
26
27         fpit = fpit->fp_next;
28     }
29
30     return 0;
31 }
```

Listing 24: vmap_page_range() in os-mm.h

4. Design Strategies Used:

The implementation employs several optimization strategies to mitigate the overhead of 5-level paging:

- **Demand Allocation:** Tables allocated only when first accessed (`get_next_level()` with `alloc=1`)
- **Sparse Page Tables:** No allocation for unused address ranges
- **4KB Alignment:** Efficient memory allocation using `posix_memalign()` for hardware compatibility
- **Present Bit Checking:** Fast detection of absent entries using `PAGING_PRESENT` bit (bit 0)
- **FIFO Page Tracking:** Efficient page replacement tracking via `fifo_pgn` linked list

In summary, moving from 2-level to 5-level paging in this Simple OS enables massive address space expansion (from 1 MB to 128 PiB) but requires more complex address parsing, slower translation (5 memory accesses vs. 2), and sophisticated memory management strategies. The implementation successfully balances these trade-offs through sparse allocation, on-demand table creation, and efficient tree traversal algorithms, making it essential for supporting modern 64-bit operating systems and large-scale applications.

4.6 Question 6

Question:What are the advantages and disadvantages of segmentation with paging?

Answer:

In the Simple OS implementation, segmentation with paging is realized through the `vm_area_struct` (VM areas) acting as logical segments, which are then divided into fixed-size pages managed by the page table hierarchy. This hybrid approach combines the benefits of both memory management techniques.

1. Implementation in Simple OS

The OS uses `vm_area_struct` (defined in `os-mm.h`) (Listing 6) as segment-like structures, each containing:

- `vm_id`: Identifier for the VM area (logical segment identifier)
- `vm_start` and `vm_end`: Boundaries defining the logical segment's virtual address range
- `sbrk`: Break pointer for dynamic memory allocation within the segment (heap growth)
- `vm_freerg_list`: Linked list of free regions within the segment for efficient reuse
- `vm_mm`: Back-pointer to the `mm_struct` for accessing page tables

Each VM area is then paged using the hierarchical page table structure (`mm_struct->pgd` in 64-bit mode, as defined in `os-mm.h`), where pages are mapped to non-contiguous physical frames. The symbol table (`symrgtbl`) tracks memory regions (variables) within VM areas, similar to how segments track logical units. The `mm_struct->mmap` field maintains a linked list of all VM areas for a process.

```
1  /*
2  * Memory management struct
3  * Represents the memory management information for a process
4  */
5  struct mm_struct
6  {
7      /* TODO: The structure of page diractory need to be justify
8       *      as your design. The single point is draft to avoid
9       *      compiler noisy only, this design need to be revised
10     */
11     #ifdef MM64
12     struct pgd_t *pgd;
13     #else
14     uint32_t *pgd;
15     #endif
16
17     struct vm_area_struct *mmap; // Linked list of memory areas
18
19     /* Currently we support a fixed number of symbol */
20     struct vm_rg_struct symrgtbl[PAGING_MAX_SYMTBL_SZ]; // Symbol region table
21
22     /* list of free page */
```

```
23     struct pgn_t *fifo_pgn; // FIFO page number list
24 };
```

Listing 25: mm_struct in os-mm.h

2. Advantages

- a. Eliminates External Fragmentation:** In `vm_map_ram()` (lines 322-337 of `mm64.c`) and `vmap_page_range()` (Listing 25), physical memory is allocated in fixed-size 4KB frames. The `framephy_struct` linked list (defined in `os-mm.h` lines 155-162) manages non-contiguous frames, eliminating external fragmentation entirely since segments are paged. The `alloc_pages_range()` function (lines 283-317) allocates frames individually and links them, allowing a VM area to span multiple non-contiguous physical frames.

```
1 // mm64.c
2 /*
3  * vm_map_ram - do the mapping all vm are to ram storage device
4  */
5 addr_t vm_map_ram(struct pcb_t *caller, addr_t astart, addr_t aend, addr_t
6     mapstart, int incpgnum, struct vm_rg_struct *ret_rg)
7 {
8     struct framephy_struct *frm_lst = NULL;
9     addr_t ret_alloc = 0;
10
11     ret_alloc = alloc_pages_range(caller, incpgnum, &frm_lst);
12
13     if (ret_alloc == -3000)
14         return -1;
15     if (ret_alloc < 0)
16         return -1;
17
18     vmap_page_range(caller, mapstart, incpgnum, frm_lst, ret_rg);
19
20     return 0;
21 }
```

Listing 26: vm_map_ram() in mm64.c

```
1 // mm64.c
2 /*
3  * alloc_pages_range - allocate req_pgnum of frame in ram
4  */
5 addr_t alloc_pages_range(struct pcb_t *caller, int req_pgnum, struct
6     framephy_struct **frm_lst)
7 {
8     int pgit;
```

```
8     addr_t fpn;
9
10    *frm_lst = NULL;
11    struct framephy_struct *tail = NULL;
12
13    for (pgit = 0; pgit < req_pgnum; pgit++)
14    {
15        if (MEMPHY_get_freefp(caller->krnl->mram, &fpn) == 0)
16        {
17            struct framephy_struct *new_node = malloc(sizeof(struct
18                framephy_struct));
19            new_node->fpn = fpn;
20            new_node->fp_next = NULL;
21            new_node->owner = caller->mm;
22
23            if (*frm_lst == NULL)
24            {
25                *frm_lst = new_node;
26                tail = new_node;
27            }
28            else
29            {
30                tail->fp_next = new_node;
31                tail = new_node;
32            }
33        }
34        else
35        {
36            return -3000;
37        }
38    }
39    return 0;
}
```

Listing 27: alloc_pages_range() in mm64.c

- b. Logical Memory Organization:** VM areas provide logical organization of memory, separating different memory regions (code, data, heap) conceptually. Functions like `get_vma_by_num()` (lines 12-31 of `mm-vm.c`) traverse the VM area list to find segments by ID, and `validate_overlap_vm_area()` (lines 76-102) ensures distinct logical regions don't overlap, allowing the OS to treat different memory areas as separate segments while using paging for physical allocation.

```
1 // mm-vm.c
2 /*get_vma_by_num - get vm area by numID
3 *@mm: memory region
```



```
4  *@vmaid: ID vm area to alloc memory region
5  *
6  */
7  struct vm_area_struct *get_vma_by_num(struct mm_struct *mm, int vmaid)
8  {
9      struct vm_area_struct *pvma = mm->mmap;
10
11     if (mm->mmap == NULL)
12         return NULL;
13
14     // Iterate through the entire list
15     while (pvma != NULL)
16     {
17         // Check for exact match
18         if (pvma->vm_id == vmaid)
19             return pvma;
20
21         pvma = pvma->vm_next;
22     }
23
24     // If we reach here, the ID was not found
25     return NULL;
26 }
```

Listing 28: get_vma_by_num() in mm-vm.c

```
1  // mm-vm.c
2  /*validate_overlap_vm_area
3  *@caller: caller
4  *@vmaid: ID vm area to alloc memory region
5  *@vmastart: vma end
6  *@vmaend: vma end
7  *
8  */
9  int validate_overlap_vm_area(struct pcb_t *caller, int vmaid, addr_t vmastart
10     , addr_t vmaend)
11  {
12     if (vmastart >= vmaend)
13     {
14         return -1; // Invalid range
15     }
16
17     struct vm_area_struct *vma = caller->mm->mmap;
18
19     // Traverse the linked list of VM Areas
```

```

19     while (vma != NULL)
20     {
21         // Skip self if we are checking against our own ID (though typically
           // we check new vs old)
22         if (vma->vm_id != vmaid)
23         {
24             // Check for overlap: OVERLAP(Start1, End1, Start2, End2)
           // Logic: (Start1 < End2) && (Start2 < End1)
25             if ((vmastart < vma->vm_end) && (vma->vm_start < vmaend))
26             {
27                 return -1; // Overlap detected
28             }
29         }
30     }
31     vma = vma->vm_next;
32 }
33
34 return 0;
35 }

```

Listing 29: validate_overlap_vm_area() in mm-vm.c

c. Efficient Handling of Large Segments: The inc_vma_limit() function (lines 110-168 of mm-vm.c) demonstrates how large VM areas can grow dynamically. Since these segments are paged, they don't require contiguous physical memory. The function:

- Calculates page-aligned size using DIV_ROUND_UP(inc_amt, PAGING_PAGESZ)
- Calls vm_map_ram() which allocates frames via alloc_pages_range() and maps them via vmap_page_range()
- Updates the sbrk pointer to track the new segment boundary

```

1           // mm-vm.c
2 /*inc_vma_limit - increase vm area limits to reserve space for new
           variable
3 *@caller: caller
4 *@vmaid: ID vm area to alloc memory region
5 *@inc_sz: increment size
6 *
7 */
8 int inc_vma_limit(struct pcb_t *caller, int vmaid, addr_t inc_sz)
9 {
10     struct vm_rg_struct *newrg = malloc(sizeof(struct vm_rg_struct));
11     int inc_amt = (int)inc_sz;
12
13     // Retrieve current VMA
14     struct vm_area_struct *cur_vma = get_vma_by_num(caller->mm, vmaid);
15     if (!cur_vma)
16     {

```

```
17         free(newrg);
18         return -1;
19     }
20
21     // Calculate number of pages needed for this increase
22     int incnumpage = DIV_ROUND_UP(inc_amt, PAGING_PAGESZ);
23
24     // Calculate new boundaries
25     struct vm_rg_struct *area = get_vm_area_node_at_brk(caller, vmaid,
26         inc_sz, inc_amt);
27     if (!area)
28     {
29         free(newrg);
30         return -1;
31     }
32
33     // Store old sbrk to define the map start
34     addr_t old_sbrk = cur_vma->sbrk;
35
36     // Validate that expanding this VMA won't collide with another VMA
37     if (validate_overlap_vm_area(caller, vmaid, area->rg_start, area->
38         rg_end) < 0)
39     {
40         free(newrg);
41         free(area);
42         return -1; /* Overlap and failed allocation */
43     }
44
45     // Call the Hardware Mapper:
46     // 1. Allocates physical frames (alloc_pages_range)
47     // 2. Updates Page Tables (vmap_page_range)
48     if (vm_map_ram(caller, area->rg_start, area->rg_end,
49         old_sbrk, incnumpage, newrg) < 0)
50     {
51         free(newrg);
52         free(area);
53         return -1; /* Map the memory to MEMRAM failed (OOM) */
54     }
55
56     // Update the VMA's Break Pointer (sbrk)
57     // This officially "commits" the expansion
58     cur_vma->sbrk += inc_sz;
59     cur_vma->vm_end += inc_sz; // Assuming vm_end tracks the used limit
60     in this simple model
```

```
58
59     // Cleanup: we used newrg in vm_map_ram via pointer,
60     // but the actual region tracking is often added to the free list
61     // or symbol table by the caller (libmem).
62     free(newrg);
63     free(area);
64
65     return 0;
66 }
```

Listing 30: inc_vma_limit() in mm-vm.c

This allows segments to grow without requiring contiguous physical memory allocation.

- d. **Reduced Internal Fragmentation:** The `vm_freerg_list` within each VM area tracks free regions, allowing efficient reuse of freed memory within segments. The `get_free_vmrg_area()` function (lines 384-437 of `libmem.c`) searches the free list for regions large enough to satisfy allocation requests, reducing internal fragmentation. The symbol table (`symrgtbl`) with `PAGING_MAX_SYMTBL_SZ = 30` entries enables tracking of variables within logical segments, and the `__alloc()` function first attempts to reuse free regions before expanding the segment.

```
1  // libmem.c
2  /* get_free_vmrg_area - get a free vm region */
3  int get_free_vmrg_area(struct pcb_t *caller, int vmaid, int size, struct
   vm_rg_struct *newrg)
4  {
5     struct vm_area_struct *cur_vma = get_vma_by_num(caller->mm, vmaid);
6     if (!cur_vma)
7         return -1;
8
9     struct vm_rg_struct *rgit = cur_vma->vm_freerg_list;
10    if (rgit == NULL)
11        return -1;
12
13    newrg->rg_start = newrg->rg_end = -1;
14
15    // Traverse free list
16    while (rgit != NULL)
17    {
18        if (rgit->rg_start + size <= rgit->rg_end)
19        {
20            /* Found fit */
21            newrg->rg_start = rgit->rg_start;
22            newrg->rg_end = rgit->rg_start + size;
23
24            /* Update free node */
```

```

25         if (rgit->rg_start + size < rgit->rg_end)
26         {
27             rgit->rg_start += size;
28         }
29         else
30         {
31             /* Perfectly matches, remove this free node */
32             struct vm_rg_struct *nextrg = rgit->rg_next;
33             if (nextrg)
34             {
35                 rgit->rg_start = nextrg->rg_start;
36                 rgit->rg_end = nextrg->rg_end;
37                 rgit->rg_next = nextrg->rg_next;
38                 free(nextrg);
39             }
40             else
41             {
42                 // This was the last node, mark as invalid/empty
43                 rgit->rg_start = rgit->rg_end = 0;
44                 rgit->rg_next = NULL;
45                 // Ideally we should remove rgit from the list head if it
46                 // 's the head
47                 // But standard linked list removal is tricky with just '
48                 // rgit' ptr
49                 // Assuming list management handles this or next alloc
50                 // handles 0-size
51             }
52         }
53         return 0;
54     }
55     rgit = rgit->rg_next;
56 }
57 return -1; // No fit found
58 }

```

Listing 31: get_free_vmrg_area() in libmem.c

- e. Support for Dynamic Growth:** The `sbrk` pointer in each VM area allows segments to grow dynamically. The `get_vm_area_node_at_brk()` function (lines 47-67 of `mm-vm.c`) creates new regions starting at the current `sbrk` position. When `get_free_vmrg_area()` fails (line 45 of `libmem.c`, inside `__alloc()` function), the system expands the heap via `inc_vma_limit()`, which calls `vm_map_ram()` to allocate and map new pages. This enables processes to extend their segments on demand without pre-allocation.

```
1 // mm-vm.c
```

```
2 struct vm_rg_struct *get_vm_area_node_at_brk(struct pcb_t *caller, int vmaid,
3     addr_t size, addr_t alignedsz)
4 {
5     struct vm_rg_struct *newrg;
6
7     // Retrieve the VMA (usually ID 0 for data/heap)
8     struct vm_area_struct *cur_vma = get_vma_by_num(caller->mm, vmaid);
9     if (!cur_vma)
10         return NULL;
11
12     newrg = malloc(sizeof(struct vm_rg_struct));
13
14     // The new region starts where the current Heap ends (sbrk)
15     newrg->rg_start = cur_vma->sbrk;
16
17     // It ends after 'size' bytes
18     newrg->rg_end = newrg->rg_start + size;
19
20     newrg->rg_next = NULL;
21
22     return newrg;
23 }
```

Listing 32: get_vm_area_node_at_brk() in mm-vm.c

```
1 // libmem.c
2 /* __alloc - allocate a region memory */
3 int __alloc(struct pcb_t *caller, int vmaid, int rgid, addr_t size, addr_t *
4     alloc_addr)
5 {
6     pthread_mutex_lock(&mmvm_lock);
7     struct vm_rg_struct rgnode;
8
9     // Check if we can reuse a free region
10    if (get_free_vmrg_area(caller, vmaid, size, &rgnode) == 0)
11    {
12        caller->mm->symrgtbl[rgid].rg_start = rgnode.rg_start;
13        caller->mm->symrgtbl[rgid].rg_end = rgnode.rg_end;
14        *alloc_addr = rgnode.rg_start;
15        pthread_mutex_unlock(&mmvm_lock);
16        return 0;
17    }
18
19    /* If get_free_vmrg_area FAILED, we must expand the heap */
20    struct vm_area_struct *cur_vma = get_vma_by_num(caller->mm, vmaid);
```

```
20     if (!cur_vma)
21     {
22         pthread_mutex_unlock(&mmvm_lock);
23         return -1;
24     }
25
26     addr_t inc_sz = PAGING_PAGE_ALIGNSZ(size);
27     // Usually aligning to page size is good practice, or use exact size if
28     // your inc_vma handles it.
29     // The provided skeleton had some ifdef logic, let's stick to standard
30     // alignment.
31
32     addr_t old_sbrk = cur_vma->sbrk;
33
34     /* SYSCALL to increase limit */
35     struct sc_regs regs;
36     regs.a1 = SYSMEM_INC_OP;
37     regs.a2 = vmaid;
38     regs.a3 = inc_sz; // Request expansion
39
40     // Note: In real syscall, we don't pass PCB, but here we simulate it via
41     // wrapper
42     // The wrapper 'syscall' takes (krnl, pid, nr, regs)
43     if (syscall(caller->krnl, caller->pid, 17, &regs) < 0)
44     {
45         pthread_mutex_unlock(&mmvm_lock);
46         return -1; // Failed to expand
47     }
48
49     /* Successful increase limit */
50     caller->mm->symrgtbl[rgid].rg_start = old_sbrk;
51     caller->mm->symrgtbl[rgid].rg_end = old_sbrk + size;
52     *alloc_addr = old_sbrk;
53
54     // The remaining space (inc_sz - size) should technically be added to
55     // free list
56     // to avoid internal fragmentation, but for this simple assignment we
57     // might skip it
58     // or add logic:
59     if (inc_sz > size)
60     {
61         struct vm_rg_struct *fragment = malloc(sizeof(struct vm_rg_struct));
62         fragment->rg_start = old_sbrk + size;
63         fragment->rg_end = old_sbrk + inc_sz;
```

```
59     fragment->rg_next = NULL;
60     enlist_vm_freerg_list(caller->mm, fragment);
61 }
62
63 pthread_mutex_unlock(&mmvm_lock);
64 return 0;
65 }
```

Listing 33: __alloc() in libmem.c

3. Disadvantages

- a. Increased Translation Complexity:** Address translation requires two conceptual steps: first identifying the VM area (segment) that contains the address, then translating the address through the page table hierarchy. In practice, the implementation uses the page table directly (via `get_page_table_entry()` in `mm64.c`, Listing 12), but VM area validation adds overhead. In 64-bit mode, address translation involves extracting indices from PGD, P4D, PUD, PMD, and PT levels (as in `get_pd_from_address()`, lines 149-157 of `mm64.c`, Listing 16), requiring 5 memory accesses. Additionally, functions like `__read()` (lines 279-290 of `libmem.c`) must first validate the address is within the correct VM area region via `get_symrg_byid()` before performing page translation.

```
1  // libmem.c
2  /* __read - read value in region memory */
3  int __read(struct pcb_t *caller, int vmaid, int rgid, addr_t offset, BYTE *
    data)
4  {
5      struct vm_rg_struct *curr = get_symrg_byid(caller->mm, rgid);
6      if (!curr)
7          return -1;
8
9      // Bounds check
10     if (curr->rg_start + offset >= curr->rg_end)
11         return -1;
12
13     return pg_getval(caller->mm, curr->rg_start + offset, data, caller);
14 }
```

Listing 34: __read() in libmem.c

- b. Memory Overhead for Management Structures:** Each process maintains a linked list of VM areas (`mm_struct->mmap`), each with its own free region list (`vm_freerg_list`). Additionally, the hierarchical page table structure (`mm_struct->pgd` in 64-bit mode) must be maintained, with each directory level consuming 4KB when allocated. The `symrgtbl` array (`PAGING_MAX_SYMTBL_SZ = 30`) adds further overhead for tracking memory regions within segments. Each VM area structure itself consumes memory, and the free region lists within each VM area add additional overhead compared to a pure paging system.

- c. Complex Overlap Validation:** The `validate_overlap_vm_area()` function in `mm-vm.c` (Listing 29) must check all VM areas to prevent overlaps, requiring traversal of the entire `mmap` linked list. This adds $O(n)$ overhead during memory allocation operations, as seen in `inc_vma_limit()` and `__alloc()` when expanding segments. The overlap check logic (lines 15-18 below) compares the new region against all existing VM areas, which becomes expensive as the number of VM areas grows.

```
1 // mm-vm.c
2 int validate_overlap_vm_area(struct pcb_t *caller, int vmaid, addr_t vmastart
  , addr_t vmaend)
3 {
4     /*
5     ...
6     */
7
8     while (vma != NULL)
9     {
10         // Skip self if we are checking against our own ID (though typically
11         // we check new vs old)
12         if (vma->vm_id != vmaid)
13         {
14             // Check for overlap: OVERLAP(Start1, End1, Start2, End2)
15             // Logic: (Start1 < End2) && (Start2 < End1)
16             if ((vmastart < vma->vm_end) && (vma->vm_start < vmaend))
17             {
18                 return -1; // Overlap detected
19             }
20             vma = vma->vm_next;
21         }
22     }
23     return 0;
24 }
```

Listing 35: Overlap check logic

- d. Page Table Management Complexity:** The `mm_struct` must maintain both VM area structures and page table structures simultaneously. In 64-bit mode, the structure contains the root PGD pointer (`mm->pgd`), and the initialization in `init_mm()` (`mm64.c`, Listing 22) shows the complexity of managing both segment-like VM areas (initializing `mmap`, `symrgtbl`) and hierarchical page tables (initializing `pgd` as `NULL`). The `vmap_page_range()` function must walk the 5-level page table hierarchy while also tracking which VM area the pages belong to, adding complexity to memory mapping operations.
- e. Swapping Complexity:** The swapping mechanism operates at the page level via `pg_getpage()` (lines 177-237 of `libmem.c`), but the VM area structure adds complexity because:

- The FIFO page replacement (`fifo_pgn` list, line 142 of `os-mm.h`) operates at the page level without direct VM area context
- When a page is swapped out, the system must ensure the page table entry is properly marked (via `pte_set_swap()`, lines 171-187 of `mm64.c`)
- The `find_victim_page()` function (lines 352-381 of `libmem.c`) selects victims from the FIFO list without considering VM area boundaries, which could lead to swapping pages from different logical segments

The `__mm_swap_page()` function (lines 33-37 of `mm-vm.c`) handles the actual swap copy, but the coordination between VM area tracking and page-level swapping adds complexity.

```
1 // libmem.c
2 /* find_victim_page - FIFO replacement policy */
3 int find_victim_page(struct mm_struct *mm, addr_t *retpgn)
4 {
5     struct pgn_t *pg = mm->fifo_pgn;
6     if (!pg)
7         return -1; // Empty list
8
9     struct pgn_t *prev = NULL;
10
11     // Traverse to the end of the list (Oldest item)
12     while (pg->pg_next)
13     {
14         prev = pg;
15         pg = pg->pg_next;
16     }
17
18     *retpgn = pg->pgn;
19
20     // Unlink
21     if (prev)
22     {
23         prev->pg_next = NULL;
24     }
25     else
26     {
27         mm->fifo_pgn = NULL; // Only one item was in list
28     }
29
30     free(pg);
31     return 0;
32 }
```

Listing 36: `find_victim_page()` in `libmem.c`

- f. Synchronization Overhead:** The `__alloc()` function uses `pthread_mutex_lock(&mmvm_lock)` to synchronize access to VM areas and page tables. This mutex protects both segment-level operations (VM area list traversal, free region management) and page-level operations (page table updates). The `__write()` function (line 312) also acquires this lock, and `__read()` performs page table operations that may require synchronization. This single mutex protects multiple data structures, potentially creating contention points that wouldn't exist in a pure paging system where only page tables need protection.
- g. Free Region List Management Overhead:** The `get_free_vmrg_area()` function (`libmem.c`) must traverse the free region list within a VM area to find suitable regions, which can be $O(n)$ in the number of free regions. When regions are freed via `__free()`, the function must create new free region nodes and insert them into the list, potentially fragmenting the free list. The `enlist_vm_freerg_list()` function adds regions to the head of the list, which is efficient but doesn't merge adjacent free regions, potentially leading to fragmentation.

In summary, the Simple OS's segmentation-with-paging approach (VM areas + page tables) provides logical memory organization and eliminates external fragmentation, but at the cost of increased translation complexity, memory overhead for management structures, $O(n)$ overlap validation, complex coordination between segment and page-level operations, and synchronization overhead. The hybrid system requires maintaining both VM area metadata and page table structures, with operations spanning both abstraction levels, making the implementation more complex than a pure paging system.

4.7 Question 7

Question: What happens if the synchronization is not handled in your Simple OS? Illustrate the problem of your simple OS (assignment outputs) by example if you have any.

Answer:

The Simple OS uses a multi-threaded architecture where multiple CPU threads (from `cpu_routine()` in `os.c`) and a loader thread (from `ld_routine()`) execute concurrently, sharing critical data structures. Without proper synchronization, severe system failures would occur.

1. Synchronization Mechanisms in Simple OS

The OS implements three main synchronization mechanisms:

- **Timer Synchronization** (`timer.c`): Uses mutexes (`event_lock`, `timer_lock`) and condition variables (`event_cond`, `timer_cond`) to coordinate time slot advancement. The `timer_routine()` waits for all devices to complete before incrementing `_time`, and `next_slot()` synchronizes CPU/loader threads with the timer.
- **Scheduler Synchronization** (`sched.c`): Uses `queue_lock` mutex to protect shared scheduler state including `mlq_ready_queue[]`, `slot_remaining[]`, `current_queue_index`, and `running_list`. Functions like `get_mlq_proc()`, `put_mlq_proc()`, and `add_mlq_proc()` are all protected by this mutex.

- **Memory Management Synchronization** (`libmem.c`): Uses `mmvm_lock` mutex to protect VM area structures and page table operations. Functions like `__alloc()`, `__write()`, and `__read()` acquire this lock before accessing memory management structures.

2. Consequences of Missing Synchronization

- a. **Ready Queue Corruption:** Without `queue_lock` protection, multiple CPU threads could simultaneously call `get_mlq_proc()` and `put_mlq_proc()`, leading to race conditions on `mlq_ready_queue[]`. This could result in:
- **Lost Processes:** A process dequeued by one CPU thread might be simultaneously dequeued by another, causing one CPU to receive `NULL` or an invalid process pointer.
 - **Duplicate Process Execution:** The same process could be dispatched to multiple CPUs simultaneously, leading to inconsistent state and data corruption.
 - **Queue Structure Corruption:** Concurrent enqueue/dequeue operations could corrupt the queue's internal structure (size, pointers), causing segmentation faults or infinite loops.
- b. **Scheduler State Inconsistency:** The scheduler maintains state variables `slot[]` and `mlq_ready_queue` that must be updated atomically. Without synchronization:
- **Incorrect Priority Scheduling:** Multiple threads reading and updating `mlq_ready_queue` simultaneously could cause the scheduler to skip priority levels or process them out of order, violating the MLQ scheduling policy.
 - **Slot Counter Corruption:** The `slot[MAX_PRIO]` counter could be incorrectly decremented by multiple threads, causing processes to receive incorrect time slices or be prematurely moved to lower priority queues.
- c. **Timer Desynchronization:** The timer mechanism in `timer_routine()` coordinates all threads by waiting for each device's `done` flag before advancing time slots. Without proper mutex protection:
- **Time Slot Skips:** If multiple threads call `next_slot()` simultaneously, the timer might advance before all threads complete their work, causing some threads to miss time slots or execute out of order.
 - **Deadlock in Timer:** The condition variable waits in `timer_routine()` (lines 18-20 below) and `next_slot()` (lines 61-63 below) could deadlock if mutexes are not properly acquired, freezing the entire system.
 - **Inconsistent Time Perception:** Different CPU threads might observe different values of `current_time()` if the timer advances while threads are reading it, leading to incorrect process scheduling decisions.

```
1 // timer.c
2 static void *timer_routine(void *args)
3 {
4     while (!timer_stop)
5     {
6         printf("Time slot %3llu\n", current_time());
7         int fsh = 0;
```

```
8      int event = 0;
9      /* Wait for all devices have done the job in current
10     * time slot */
11     struct timer_id_container_t *temp;
12     for (temp = dev_list; temp != NULL; temp = temp->next)
13     {
14         pthread_mutex_lock(&temp->id.event_lock);
15         while (!temp->id.done && !temp->id.fsh)
16         {
17             pthread_cond_wait(
18                 &temp->id.event_cond,
19                 &temp->id.event_lock);
20         }
21         if (temp->id.fsh)
22         {
23             fsh++;
24         }
25         event++;
26         pthread_mutex_unlock(&temp->id.event_lock);
27     }
28
29     /* Increase the time slot */
30     _time++;
31
32     /* Let devices continue their job */
33     for (temp = dev_list; temp != NULL; temp = temp->next)
34     {
35         pthread_mutex_lock(&temp->id.timer_lock);
36         temp->id.done = 0;
37         pthread_cond_signal(&temp->id.timer_cond);
38         pthread_mutex_unlock(&temp->id.timer_lock);
39     }
40     if (fsh == event)
41     {
42         break;
43     }
44 }
45 pthread_exit(args);
46 }
47
48 void next_slot(struct timer_id_t *timer_id)
49 {
50     /* Tell to timer that we have done our job in current slot */
51     pthread_mutex_lock(&timer_id->event_lock);
```

```
52     timer_id->done = 1;
53     pthread_cond_signal(&timer_id->event_cond);
54     pthread_mutex_unlock(&timer_id->event_lock);
55
56     /* Wait for going to next slot */
57     pthread_mutex_lock(&timer_id->timer_lock);
58     while (timer_id->done)
59     {
60         pthread_cond_wait(
61             &timer_id->timer_cond,
62             &timer_id->timer_lock);
63     }
64     pthread_mutex_unlock(&timer_id->timer_lock);
65 }
```

Listing 37: timer_routine() and next_slot() in timer.c

- d. **Memory Management Corruption:** Without `mmvm_lock` protection, concurrent memory operations could corrupt critical structures:
- **VM Area List Corruption:** Multiple threads accessing `mm_struct->mmap` (the VM area linked list) simultaneously could corrupt list pointers, causing processes to access wrong memory regions or crash when traversing the list.
 - **Page Table Race Conditions:** Concurrent calls to `pte_set_fpn()` and `pte_get_entry()` could result in reading partially updated page table entries, causing incorrect address translations and memory access violations.
 - **Symbol Table Corruption:** The `symrgtbl[]` array could be corrupted if multiple threads allocate memory regions simultaneously, leading to processes accessing incorrect memory addresses or overwriting each other's data.
- e. **Process State Corruption:** In `cpu_routine()`, each CPU thread maintains local state (`proc`, `time_left`) but accesses shared process control blocks. Without synchronization:
- **Concurrent Process Execution:** The same `pcb_t` structure could be modified by multiple CPU threads simultaneously, corrupting the program counter (`pc`), registers, or memory mappings.
 - **Inconsistent Process Termination:** A process might be freed by one CPU thread while another is still executing it, leading to use-after-free errors and system crashes.
- f. **Loader-CPU Race Conditions:** The `ld_routine()` adds processes to the ready queue via `add_proc()` while CPU threads are simultaneously removing processes via `get_proc()`. Without synchronization:
- **Process Addition During Dequeue:** A process might be added to the queue while another thread is in the middle of dequeuing, causing the new process to be lost or the dequeue operation to fail.
 - **Incomplete Process Initialization:** A CPU thread might dispatch a process before the loader

finishes initializing its memory structures (`mm_struct`), leading to null pointer dereferences or uninitialized memory access.

3. System-Wide Failures

Without synchronization, the Simple OS would experience cascading failures:

- **Non-Deterministic Behavior:** The same input could produce different outputs on each run due to race conditions, making debugging and testing extremely difficult.
- **System Crashes:** Segmentation faults from corrupted pointers, null pointer dereferences, or invalid memory accesses would cause the entire OS simulation to crash.
- **Resource Leaks:** Lost processes or corrupted queue structures could prevent proper cleanup, leading to memory leaks and resource exhaustion.
- **Incorrect Execution Results:** Even if the system doesn't crash, processes might execute incorrectly due to corrupted state, producing wrong outputs that are difficult to detect.

4. Experimental Demonstration: Removing Timer Synchronization

To demonstrate the consequences of missing synchronization, we removed all mutex protection from `timer_routine()` (`timer.c`) and `next_slot()` functions. This eliminates the critical synchronization mechanism that coordinates time slot advancement across all CPU threads and the loader thread.

a. Experimental Setup

The mutexes `event_lock` and `timer_lock` (defined per device in `struct timer_id_t`), along with their associated condition variables (`event_cond` and `timer_cond`), were removed from the timer synchronization code. Specifically:

- In `timer_routine()`: Removed mutex locks at lines 32, 44, 53, 56 that protect the `done` flag checks and updates

```
1 // timer.c
2 static void *timer_routine(void *args)
3 {
4     /*...*/
5     for (temp = dev_list; temp != NULL; temp = temp->next)
6     {
7 //line 32  pthread_mutex_lock(&temp->id.event_lock);
8         while (!temp->id.done && !temp->id.fsh)
9         {
10             pthread_cond_wait(
11                 &temp->id.event_cond,
12                 &temp->id.event_lock);
13         }
14         if (temp->id.fsh)
15         {
16             fsh++;
17         }
18     }
```

```

18         event++;
19 //line 44  pthread_mutex_unlock(&temp->id.event_lock);
20     }
21     /*...*/
22     for (temp = dev_list; temp != NULL; temp = temp->next)
23     {
24 //line 53  pthread_mutex_lock(&temp->id.timer_lock);
25         temp->id.done = 0;
26         pthread_cond_signal(&temp->id.timer_cond);
27 //line 56  pthread_mutex_unlock(&temp->id.timer_lock);
28     }
29     /*...*/
30 }

```

Listing 38: Removed muxtex timer_routine()

- In next_slot(): Removed mutex locks at lines 69, 72, 75, 82 that coordinate between threads and the timer

```

1 // timer.c
2 void next_slot(struct timer_id_t *timer_id)
3 {
4     /* Tell to timer that we have done our job in current slot */
5 //line 69  pthread_mutex_lock(&timer_id->event_lock);
6         timer_id->done = 1;
7         pthread_cond_signal(&timer_id->event_cond);
8 //line 72  pthread_mutex_unlock(&timer_id->event_lock);
9
10    /* Wait for going to next slot */
11 //line 75  pthread_mutex_lock(&timer_id->timer_lock);
12    while (timer_id->done)
13    {
14        pthread_cond_wait(
15            &timer_id->timer_cond,
16            &timer_id->timer_lock);
17    }
18 //line 82  pthread_mutex_unlock(&timer_id->timer_lock);
19 }

```

Listing 39: Removed muxtex next_slot()

This allows multiple threads (CPU threads from `cpu_routine()` in `os.c` and the loader thread from `ld_routine()`) to access and modify the timer state (`_time` of `timer.c`) and device `done` flags concurrently without any protection.

b. Observed Output

The following output was captured from a test run with timer synchronization removed:


```
1 Time slot    0
2 ld_routine
3     Loaded a process at input/proc/p0s, PID: 1 PRI0: 0
4
5 Time slot    1
6     CPU 0: Dispatched process  1
7     Loaded a process at input/proc/p1s, PID: 2 PRI0: 15
8     CPU 1: Dispatched process  2
9
10 Time slot   2
11 liballoc:152
12 --- Page Table Dump ---
13 print_pgtbl:
14     PGD=0x7fbab8001000 P4D=0x7fbab8003000 PUD=0x7fbab8005000 PMD=0x7fbab8007000
15     PT=0x7fbab8009000
16     0000000000000000: 0000000000000001
17
18 Time slot   3
19 liballoc:152
20 --- Page Table Dump ---
21 print_pgtbl:
22     PGD=0x7fbab8001000 P4D=0x7fbab8003000 PUD=0x7fbab8005000 PMD=0x7fbab8007000
23     PT=0x7fbab8009000
24     0000000000000000: 0000000000000001
25
26 Time slot   4
27 libfree:168
28 --- Page Table Dump ---
29 print_pgtbl:
30     PGD=0x7fbab8001000 P4D=0x7fbab8003000 PUD=0x7fbab8005000 PMD=0x7fbab8007000
31     PT=0x7fbab8009000
32     0000000000000000: 0000000000000001
33
34 Time slot   5
35 liballoc:152
36 --- Page Table Dump ---
37 print_pgtbl:
38     PGD=0x7fbab8001000 P4D=0x7fbab8003000 PUD=0x7fbab8005000 PMD=0x7fbab8007000
39     PT=0x7fbab8009000
40     0000000000000000: 0000000000000001
41
42 Time slot   6
43 libwrite:343
```

```
41 --- Page Table Dump ---
42 print_pttbl:
43   PGD=0x7fbab8001000 P4D=0x7fbab8003000 PUD=0x7fbab8005000 PMD=0x7fbab8007000
   PT=0x7fbab8009000
44 0000000000000000: 0000000000000001
45   Loaded a process at input/proc/p0s, PID: 4 PRI0: 0
46   CPU 0: Put process 1 to run queue
47   CPU 1: Put process 2 to run queue
48   CPU 1: Dispatched process 3
49   CPU 0: Dispatched process 4
50
51 Time slot 7
52
53 Time slot 8
54 liballoc:152
55 --- Page Table Dump ---
56 print_pttbl:
57   PGD=0x7fbab800b000 P4D=0x7fbab800d000 PUD=0x7fbab800f000 PMD=0x7fbab8011000
   PT=0x7fbab8013000
58 0000000000000000: 0000000000001001
59
60 Time slot 9
61 liballoc:152
62 --- Page Table Dump ---
63 print_pttbl:
64   PGD=0x7fbab800b000 P4D=0x7fbab800d000 PUD=0x7fbab800f000 PMD=0x7fbab8011000
   PT=0x7fbab8013000
65 0000000000000000: 0000000000001001
66
67 Time slot 10
68 libfree:168
69 --- Page Table Dump ---
70 print_pttbl:
71   PGD=0x7fbab800b000 P4D=0x7fbab800d000 PUD=0x7fbab800f000 PMD=0x7fbab8011000
   PT=0x7fbab8013000
72 0000000000000000: 0000000000001001
73
74 Time slot 11
75 liballoc:152
76 --- Page Table Dump ---
77 print_pttbl:
78   PGD=0x7fbab800b000 P4D=0x7fbab800d000 PUD=0x7fbab800f000 PMD=0x7fbab8013000
79 0000000000000000: 0000000000001001
80
```

```
81 Time slot    12
82 libwrite:343
83 --- Page Table Dump ---
84 print_pgtbl:
85   PGD=0x7fbab800b000 P4D=0x7fbab800d000 PUD=0x7fbab800f000 PMD=0x7fbab8011000
      PT=0x7fbab8013000
86 0000000000000000: 00000000000001001
87
88 Time slot    13
89   CPU 1: Put process 3 to run queue
90   CPU 1: Dispatched process 1
91 libread:304
92   CPU 0: Put process 4 to run queue
93   CPU 0: Dispatched process 3
94
95 Time slot    14
96
97 Time slot    15
98
99 Time slot    16
100
101 Time slot    17
102   CPU 0: Processed 3 has finished
103   CPU 0: Dispatched process 4
104 libread:304
105
106 Time slot    18
107
108 Time slot    19
109   CPU 1: Put process 1 to run queue
110   CPU 1: Dispatched process 1
111 libfree:168
112 --- Page Table Dump ---
113 print_pgtbl:
114   PGD=0x7fbab8001000 P4D=0x7fbab8003000 PUD=0x7fbab8005000 PMD=0x7fbab8007000
      PT=0x7fbab8009000
115 0000000000000000: 0000000000000001
116
117 Time slot    20
118
119 Time slot    21
120   CPU 1: Processed 1 has finished
121   CPU 1: Dispatched process 2
122
```

```
123 Time slot    22
124
125 Time slot    23
126     CPU 0: Put process  4 to run queue
127     CPU 0: Dispatched process  4
128 libfree:168
129 --- Page Table Dump ---
130 print_pgtbl:
131     PGD=0x7fbab800b000 P4D=0x7fbab800d000 PUD=0x7fbab800f000 PMD=0x7fbab8011000
        PT=0x7fbab8013000
132 0000000000000000: 0000000000001001
133
134 Time slot    24
135
136 Time slot    25
137     CPU 0: Processed  4 has finished
138     CPU 1: Processed  2 has finished
139     CPU 0 stopped
140     CPU 1 stopped
```

c. Analysis of the Results

While this particular run completed successfully, several critical observations highlight the dangers of missing synchronization:

- **Concurrent Operations in Same Time Slot:** At time slot 6, we observe multiple events occurring simultaneously within the same time slot: `libwrite:343`, process loading (PID 4), CPU 0 putting process 1 to run queue, CPU 1 putting process 2 to run queue, and both CPUs dispatching new processes (3 and 4). Without timer synchronization, the loader thread, CPU 0, and CPU 1 are all executing concurrently. The `next_slot()` function (called from `cpu_routine()` at lines 66, 99, 112 and from `ld_routine()` at lines 140, 155 in `os.c`) should coordinate with `timer_routine()` to ensure only one time slot is active at a time. Without mutexes, threads can read and modify the `done` flag and `_time` variable concurrently, causing operations from different logical time slots to interleave.
- **Potential Page Table Race Conditions:** While the page table dumps appear consistent in this run, without proper synchronization, concurrent access to page table structures could lead to race conditions. Multiple threads could simultaneously call `pte_set_fpn()` or `pte_get_entry()` (in `mm64.c`), which traverse the 5-level page table hierarchy via `get_page_table_entry()`. Without mutex protection, one thread could be updating a page table entry (via `get_next_level()` allocating new directory levels) while another thread reads it, potentially observing partially updated pointers or inconsistent states. The fact that this run shows consistent page tables does not guarantee correctness—race conditions are non-deterministic and may manifest differently in subsequent runs.

- **Concurrent Memory Operations:** At time slot 13, we see both CPUs performing operations simultaneously: CPU 1 putting process 3 to run queue and dispatching process 1, CPU 1 performing `libread:304`, CPU 0 putting process 4 to run queue and dispatching process 3. CPU 1 is reading memory while CPU 0 is dispatching a process. Without proper synchronization, these operations could interfere with each other, especially if they access shared scheduler structures or memory management data.
- **Non-Deterministic Execution Order:** The output shows that operations from different threads can appear in any order within a time slot. For example, at time slot 6, the process loading, CPU queue operations, and dispatches all occur in an unpredictable sequence. This non-determinism makes it impossible to reproduce bugs consistently, as the exact interleaving depends on thread scheduling by the operating system.
- **Time Slot Coordination Failure:** The `timer_routine()` function (lines 20-64 of `timer.c`) is designed to:
 - i. Wait for all devices to set their `done` flag (lines 30-45)
 - ii. Increment `_time` (line 48)
 - iii. Signal all devices to continue (lines 51-57)

Without mutex protection, multiple threads can call `next_slot()` simultaneously, and the timer thread can advance `_time` while other threads are still reading it or setting their `done` flags. This breaks the atomicity of time slot transitions, causing threads to operate on inconsistent time values.

- **The Illusion of Correctness:** The fact that this run completed successfully demonstrates a critical characteristic of race conditions: they are **non-deterministic**. The absence of visible errors in one run does not guarantee correctness. While this particular run shows consistent page table structures and no obvious data corruption, the concurrent operations observed (especially at time slots 6 and 13) indicate that threads are executing without proper coordination. Subsequent runs with the same input might:
 - Crash due to corrupted pointers or invalid memory accesses
 - Produce incorrect results due to data corruption
 - Deadlock if threads wait on condition variables without proper mutex protection (the condition variable waits in `timer_routine()` and `next_slot()` require mutexes to function correctly)
 - Exhibit different execution orders, making debugging and testing impossible

d. Why This Demonstrates the Problem

The removal of timer synchronization mutexes creates a **time-of-check to time-of-use (TOC-TOU)** vulnerability and breaks the atomicity of time slot transitions. The synchronization protocol in `timer.c` works as follows:

- i. Each device thread calls `next_slot()` which sets `done = 1` and signals the timer (lines 69-72)
- ii. The device then waits for the timer to advance (lines 75-82)
- iii. The timer waits for all devices to set `done = 1` (lines 30-45)

- iv. The timer increments `_time` and resets all `done` flags (lines 48, 54-56)
- v. The timer signals all devices to continue (line 55)

Without mutex protection, this protocol breaks:

- Thread A might read `_time = 5` and set `done = 1`
- Thread B might read `_time = 5` and also set `done = 1`
- The timer thread might increment `_time` to 6 while Thread A is still in the middle of its operation
- Thread A completes its operation thinking it's in time slot 5, but the timer has already advanced
- Multiple threads can observe different values of `_time` simultaneously

The concurrent operations observed at time slots 6 and 13 provide concrete evidence of this race condition: multiple threads are executing operations simultaneously without proper coordination, breaking the atomicity of time slot transitions.

In summary, synchronization is essential for the Simple OS's correctness. The mutex-protected critical sections in `timer.c` (protecting `_time` and `done` flags), `sched.c` (protecting queue operations), and `libmem.c` (protecting memory management structures) ensure that shared data structures are accessed atomically, preventing race conditions, data corruption, and system failures. The experimental removal of timer synchronization demonstrates that even when a run appears successful, the underlying race conditions create non-deterministic behavior. The concurrent operations observed in the output show that threads are not properly coordinated, which can lead to subtle bugs, data corruption, and unpredictable system behavior in subsequent runs.

5 Overall

In this assignment, we undertook the design and implementation of the core components for a simple Operating System simulation, focusing on a Symmetric Multiprocessing (SMP) environment. The project integrated three critical modules: a Slot-based Scheduler, a Hybrid Memory Management system, and a Synchronization mechanism.

For process management, we implemented a Slot-based Multi-Level Queue (MLQ) scheduling algorithm. Unlike standard priority scheduling which is prone to starvation, our design introduces a dynamic slot mechanism ($slot = MAX_PRIO - prio$) combined with Round-Robin execution within each queue. This approach is designed to balance honoring process priority with maintaining fairness, aiming to allow lower-priority tasks to eventually receive CPU resources even under heavy load. The implementation addresses the logic of context switching and load balancing across multiple simulated CPUs.

Regarding memory management, we constructed a hybrid architecture that combines Segmentation with 5-Level Paging. A focal point of this implementation is the simulation of the hardware page walk mechanism. We programmed the traversal logic that iteratively resolves a 64-bit virtual address through the five hierarchical levels: $PGD \rightarrow P4D \rightarrow PUD \rightarrow PMD \rightarrow PT$. By utilizing bit-masking and shifting operations to extract directory indices at each level, our system replicates the translation process of modern MMUs. This mechanism is intended to enable the OS to handle sparse memory allocation within a massive 128 PiB address space while establishing isolation between user and kernel modes.

Furthermore, we addressed the challenge of concurrency in a multi-threaded architecture. We identified potential race conditions that could lead to data corruption or system instability when multiple CPUs and the loader access shared resources simultaneously. By applying synchronization mechanisms, specifically Mutex locks on the ready queue, memory structures, and the global timer, we sought to preserve the atomicity of critical operations and promote the deterministic behavior of the system.

In conclusion, this assignment has provided a comprehensive hands-on experience in building OS kernel components. It bridged the gap between theoretical concepts, such as virtual memory translation, process scheduling, thread safety and practical system programming. The insights gained from handling the complexity of the 5-level page walking algorithm and SMP synchronization serve as a foundation for understanding the architecture of modern, high-performance operating systems.

6 Source Code

1. **Scheduler:** [Link GitHub](#)
2. **Memory Management:** [Link GitHub](#)