

## CS 423 Midterm (Close Book; No Cheat Sheet)

**10/23/2025 (Thursday), In Class**

### Tips:

- **The exam duration is 1 hour.**
- **You are expected to do the exam independently.** No discussion is allowed.
- **You are not allowed to use unauthorized materials.** Textbook, class notes, cheat sheets, and electronic devices are prohibited.
- **Please write clearly and legibly.** If we can't understand your writing, then we won't give points. We won't guess anything beyond what you wrote.
- Give **short, concise answers** rather than long, vague ones (we grade by correctness, not by length).
- If you find that you have to make assumptions to solve problems, write them down (e.g., "I assume this is an x86 architecture" or "I assume a Linux kernel version >2.6.3"). Most problems are kept at a high level and do not need specific assumptions.

Sections (20 Points + 2 Bonus Points)	Points
1. Kernel Interface (5 points)	
2. Memory Management (7 points)	
3. MP-1 Continued (5 points)	
4. MP-2 Continued (5 points)	

**Name:**\_\_\_\_Peizhe Liu\_\_\_\_ **NetID:**\_\_\_\_peizhel2\_\_\_\_

## 1. Kernel Interface (5 points)

We learned in the class that `strace` is the tool for tracing system calls from a user application. Cathy wants to trace system calls for `ls`, so she runs:

```
strace ls t > /dev/null
```

And below are partial traces given by `strace`:

```
execve("/usr/bin/ls", ["ls", "t"], 0x7ffe9be0f1f8 /* 82 vars */) = 0
access("/etc/ld.so.preload", R_OK)      = -1 ENOENT (No such file or
directory)
openat(AT_FDCWD, "/etc/ld.so.cache", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st_mode=S_IFREG|0644, st_size=211975, ...}) = 0
mmap(NULL, 211975, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7f825c07a000
close(3)                                = 0
write(2, "cannot access 't'", 17cannot access 't')      = 17
write(2, ": No such file or directory", 27: No such file or directory) = 27
write(2, "\n", 1)                                = 1
exit_group(2)                             = ?
```

(1) Please explain what `mmap` does in line 6. (1 point)

This `mmap` call maps the `ld.so.cache` library into the memory.

(Dynamic linker cache)

Rubric: mention maps file into memory, mention context (FD 3, or file path), +1 point.

Mention wrong `mmap` purpose, did not mention context at all, +0 point.

(2) There is a `write` system call in line 8. Please explain what the first argument (the number 2) stands for, and explain how **this** `write` system call works. (2 points)

The first argument is the FD number of `stderr`.

This `write` call will first trap into the kernel. Via the kernel VFS layer, the terminal driver (tty) `write` function will actually handle it, printing the error message.

Rubric: mention anything like console, terminal, `std`, `stderr`, +1 point.

Briefly mention how `write` call works, mention trap into kernel, mention FD table, syscall table, VFS, TTY driver, or anything reasonable, +1 point.

Mention what `write` call does, +0 point. We ask "How," not "What."

- (3) In x86-64, the `syscall` instruction is a serializing instruction, which means that it cannot execute until all the older instructions are completed, and that it prevents the execution of all the younger instructions until it completes. Please explain why this is necessary. (2 points)

Syscall involves a context switch between kernel and user space. To prevent inconsistency between kernel and user space, the CPU flushes its pipeline.

Rubric: mention context switch, kernel trap, different runlevel +2 points

Mention data race, or pure "inconsistency" +0 point. The key here is context switch.

## 2. Memory Management (7 points)

We learned the mechanism of virtual memory management in the class based on the x86-64 semantics. As we know, memory capacity has been growing and terabyte-scale memory is already available on commodity servers. So, four-level page tables are no longer sufficient. As a result, Linux recently supported five-level page tables.

Let's think of a five-level page table, where each page-table page is 4KB, each page table entry is 8 bytes, and huge pages are 2MB and 1GB. We use L1 to refer to the last level of the page table, and L5 to refer to the highest level.

- (1) How large virtual memory can this five-level radix-tree page table support? (1 point)

4 KB page size has 12 bits of offset bits.

$4\text{ KB} / 8\text{ B} = 512$  ( $2^9$ ) entries per PT page. 5 levels yield  $5 * 9 = 45$  bits of index bits.

So we have  $12 + 45 = 57$  bit virtual addresses.

Virtual address space:  $2^{57}\text{ B} = 128\text{ PB}$ .

Rubric:  $2^{57}\text{ B}$  or  $128\text{ PB}$  +1 point. Only grade final answer, no partial points.

- (2) We learned in the class that walking over a multi-level page table is slow. Apparently, a five-level page table is slower than a four-level page table. A proposal from the class is to merge the intermediate levels to reduce the height of the page table tree. Let's say we merge L2 and L3 and turn the page table tree back into four levels: L5, L4, merged(L3, L2), and L1. In this case, what is the minimal amount of memory a process will need to use for its page table? (2 points)

merged (L3, L2) has  $2^{18}$  entries per PT page.

For 8B entry, the PT page size would be  $2^{18}\text{ Entries} * 8\text{ B Entry Size} = 2^{21}\text{ B} = 2\text{ MB}$ .

Total minimum PT size:  $4\text{ KB} * 3 + 2\text{ MB}$ .

Rubric:  $4\text{ KB} * 3 + 2\text{ MB}$ , +2 points. Only grade final answer, no partial points.

- (3) Besides improving the hardware (e.g., TLB and PWC), can you give an idea to reduce off-TLB memory translation overhead, which only needs to change software? Please explain your idea, why does it help? (2 points)

Answer not limited to: hugepages, prefault, madvise, numactl.

Rubric: mention hugepage, reduce tree level, flatten page table, improved memory allocator, improved software locality, or anything reasonable, +2 points.

Mention anything HW-required like improved cache, improved TLB, hashing-based paging, +0 point. LRU is an algorithm used in TLB.

- (4) In this five-level page table, can you describe what is the most expensive page fault and what is the overhead to address the page fault? Please give a detailed explanation (2 points)

Answer not limited to: page swap while thrashing, remote page fault.

Rubric: mention page swap, page on disk, or similar reasonable, +2 points.

Mention only L5 missing page, +0 point. Missing pages are not necessarily the most expensive. We introduced page swapping in class.

### 3. MP-1 Continued (5 points)

In MP-1, we have implemented a kernel module that measures the user space CPU time of a list of registered processes.

Cathy uses a single global lock to protect her MP-1 process list. When updating the CPU time for each process, Cathy first holds the global lock, and then iterates through the process list. While iterating, she updates the CPU time for each process. Finally, she releases the global lock.

- (1) Peizhe thinks the critical section is too long. He argues that Cathy can unlock while updating the CPU time for each process, and then lock again to continue the iteration. Do you agree with his argument? If yes, explain how this can improve the performance. If not, explain why this does not work. (1 point)

No, if the node being updated is **deleted/deallocated**, we have use after free problem.

Rubric: mention delete/deallocate, +1 point.

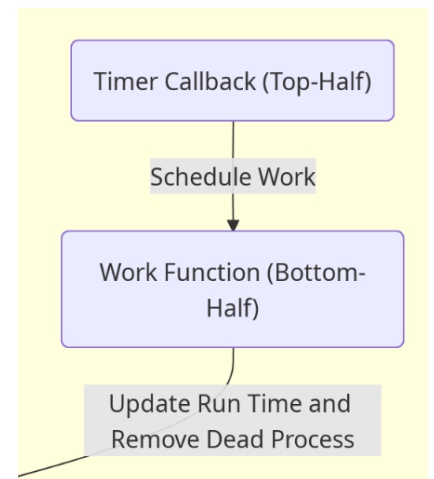
Mention "data race", "WAW" or "inaccurate cup time" +0 point. Recall we use a timer for every 5s. Only mentioning data race is not sufficient, you need to tell us what race.

Remember that we are required to use the 2-halves approach in MP-1, which is shown on the right. Timer callback, which is the first half, will schedule the work function, which is the second half.

- (2) Why do we use such a 2-halves approach, instead of implementing a monolithic timer callback function? Please explain. (1 point)

Keep the softirq section short as softirq must finish quickly.

Rubric: mention each +0.5 points. Paraphrase is ok.



Additionally, timer callback is based on softirq, which makes it an interrupt context. As we emphasized in the MP-1 documentation, interrupt context has many limitations.

- (3) We can use spinlock in timer callback, but we cannot use mutex. Please explain the reason. (1 point)

Mutex can sleep, we can't sleep/yield in the interrupt context. Spinlock does not sleep.

OR

Mutex can deadlock in an interrupt context.

Rubric: mention either above +1 point.

When grading MP-1 submissions, Cathy had an idea of exploring a new way to test the result by calling `tail -f /proc/mp1/status`. But to her surprise, this crashed many submissions' kernels even when `cat /proc/mp1/status` worked as expected. She explored a bit by checking the system call traces:

```
strace cat /proc/mp1/status > /dev/null
```

```
execve("/bin/cat", ["cat", "/proc/mp1/status"], ...) = 0
# .. traces omitted
fstat(1, {...}) = 0
openat(AT_FDCWD, "/proc/mp1/status", O_RDONLY) = 3
fstat(3, {st_mode=S_IFREG|0666, st_size=0, ...}) = 0
read(3, "", 262144) = 0
```

```
strace tail /proc/mp1/status > /dev/null
```

```
execve("/bin/tail", ["tail", "/proc/mp1/status"], ...) = 0
# .. traces omitted
openat(AT_FDCWD, "/proc/mp1/status", O_RDONLY) = 3
fstat(3, {st_mode=S_IFREG|0666, st_size=0, ...}) = 0
lseek(3, 0, SEEK_CUR[ 16.270596] BUG: kernel NULL pointer dereference,
address: 0000000000000000
# .. panic message omitted
[ 16.273918] Kernel panic - not syncing: Fatal exception
[ 16.273918] Kernel Offset: disabled
[ 16.273918] ---[ end Kernel panic - not syncing: Fatal exception ]---
```

Oh, no! The kernel crashed during the strace run. Cathy is very curious about the root cause, so she brought up the Linux kernel source, and looked for the `proc_ops` interface as a reference:

```
struct proc_ops {
    int (*proc_open)(struct inode *, struct file *);
    ssize_t (*proc_read)(struct file *, char __user *, size_t, loff_t *);
    ssize_t (*proc_read_iter)(struct kiocb *, struct iov_iter *);
    ssize_t (*proc_write)(struct file *, const char __user *, size_t, loff_t *);
    loff_t (*proc_lseek)(struct file *, loff_t, int);
    int (*proc_release)(struct inode *, struct file *);
    // .. omitted
};
```

- (4) Please help Cathy identify the possible cause for the kernel panic. You can directly mark on the logs/code above and briefly explain the cause. (2 points)

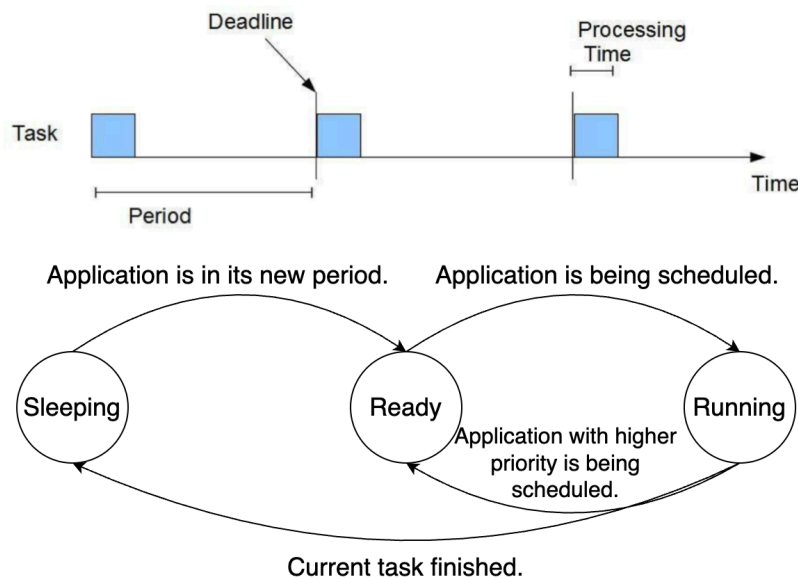
Marked on the code. `proc_lseek` function pointer is NULL, implementation is not complete.

Rubric: mention `lseek` / `proc_lseek`, +1 point.

Mention not-implemented / buggy implementation, +1 point.

#### 4. MP-2 Continued (5 points)

In MP-2, we have implemented a Rate Monotonic Scheduler (RMS), where each task is modeled by its period and computation time, and have three states: Sleeping, Ready, and Running. The task model, state machine, and state transition diagram is shown as follows:



Remember that each task must call `yield` after it finishes in the current period. In this section, you will work on the yield handler with your hardworking TA Peizhe, and explore the early/late yield problem. Peizhe's current implementation is provided below.

```
static int mp2_process_yield(int pid)
{
    struct mp2_pcb *process;
    mutex_lock(&mp2_pcb_list_lock);
    if (!(process = mp2_process_lookup_locked(pid))) { // Look up the PCB
        mutex_unlock(&mp2_pcb_list_lock);
        return -ESRCH;
    }
    mp2_process_sleep_locked(process); // Put the task to sleep and set state
    mod_timer(&process->timer, jiffies + msecs_to_jiffies(process->period));
    mutex_unlock(&mp2_pcb_list_lock);
    wake_up_process(mp2_dispatcher_kthread);
    return 0;
}
```

- (1) Peizhe noticed that his RMS was not behaving properly. **Initially all tasks run fine, but they will eventually miss deadlines, breaking the task model.** How can it break real-world RT applications? Please provide an example. (1 point)

For example: surveillance cameras can have bad frames. (1 point)

**OR**

For example: missing deadlines leads to frame losses in surveillance cameras. (1 point)

Rubric: You can either give an example of broken real-world application, or explain how it can break real-world application. Either +1 point. Not mentioning anything related to RT application +0 point.

- (2) Considering the problem symptom and his implementation, please help him identify the problem. You can directly mark on the code at the previous page. (1 point)

Marked on code. The timer was set incorrectly, resulting in later and later period starts.

Rubric: mark exactly as above +1 point, or mark the whole mod\_timer with explanation of wrong time +1 point. Only mark mod\_timer +0.5 point.

Mark wrong code or more than mod\_timer, +0 point. Mark correctly but explain incorrectly, +0 point.

As we already emphasized in the MP-2 documentation and walkthrough, your user application must register its computation time accurately. We now define the **early/late yield problem**: the user application calls yield earlier or later than its registered computation time.

- (3) What are the consequences of the early/late yield problem in RMS? Please provide one example for each case (early/late). (2 points)

**Early:** wasting CPU, lower utilization, bad admission control due to less opportunity (assume yield after computation done), or reasonable real world example.

**OR**

**Early:** task may not finish task (assume yield before computation done), or example.

**Late:** missed deadline, starved other program, bad admission control due to wrong upper bound, or reasonable real world example.

Rubric: mention either early consequence, +1 point. Mention either late consequence, +1 point. Technically we assume we yield after each computation, but we decided to give.

- (4) Thank you, now we are able to detect the problem. Peizhe wants to resolve this problem by updating the inaccurate computation time in PCB. Do you agree with his solution? If yes, explain how this can avoid consequences above. If not, explain why this does not work. (1 points)

No, this does not work. This will break the admission control.

Rubric: mention admission control, or other program missing deadline, or starve, +1 point.

Answer NO without correct explanation +0 point. Answer YES +0 point.