

# COMP 3511 Operating System (Spring 2023)

## Midterm Exam

Date: 03-Apr-2023 (Mon)

Time: 19:00 – 21:00 (2 hours)

Name (Write the name printed on your Student ID card)	<b>Solution (v7) - After grading (added partial marking schemes)</b>
Student ID	
ITSC email	

Exam format and rules:

- It is an open-book, open-notes exam (Reference: Chapter 1)
- Answer all questions within the space provided on the examination paper.
- The exam booklet is single-sided. You may use the back of the pages for your rough work.
- Please read each question very carefully and answer the question clearly to the point.
- Make sure that your answers are neatly written, legible, and readable.
- An electronic calculator is allowed, but not other electronic devices (e.g., mobile phones, tablets, ... are not allowed)

## Part I. Multiple Choices [25 \* 1 points]

Write down your answers in the boxes below:

MC1	MC2	MC3	MC4	MC5	MC6	MC7	MC8	MC9	MC10
D	B	A	C	C	D	A	A	D	B

MC11	MC12	MC13	MC14	MC15	MC16	MC17	MC18	MC19	MC20
A	B	C	C	B	A	D	B	C	Cancelled

MC21	MC22	MC23	MC24	MC25
D	A	B	C	D

### [Introduction]

MC1. Which of the following components is not a part of the Central Processing Unit in the Von Neumann architecture?

- A. Control unit
- B. Registers
- C. Arithmetic Logic Unit
- D. Main memory

Answer: D

Unlike registers, the main memory locates outside the CPU. See 1.26.

MC2. Which of the following is true about the storage system?

- A. Electrical storage is generally larger and less expensive per byte than mechanical storage.
- B. Temporal locality in caching exploited the tendency of programs to access the same data item over and over again within a small time duration.
- C. Non-Uniform Memory Access is a major approach in system design that bypasses the CPU in the data transfer between the I/O device and memory.
- D. In the caching process, the recently accessed data items tend to be copied from smaller to larger storage for sufficient storage capacity.

Answer: B

A. Electrical storage is typically costly, smaller, more reliable, and faster than mechanical storage. See 1.35.

C. The description is about DMA rather than NUMA. See 1.40 and 1.46.

D. Cache is usually much smaller than storage. See 1.36

MC3. Which of the following statement is false about multiprocessor systems?

- A. A multiprocessor system significantly saves power consumption compared to a

single-processor system during the same time period.

B. In a typical symmetric multiprocessor system, the processors usually have their own local cache while sharing the system bus.

C. Contention for shared resources lowers the expected performance gain from additional processors in multiprocessor systems.

D. Multiprocessor system can be implemented with a single chip containing multiple computing cores.

**Answer: A**

**Multiprocessor systems are economy of scale. However, it does not mean that they consume less power. Generally, more processors lead to more power consumption.**

## **[Operating System Structures]**

MC4. Which of the following is false about operating system services?

A. OS services are a set of functions that are helpful to the user or ensure the system's efficient operation via resource management.

B. GUI interface is not a mandatory part of OS services.

C. To execute a C program, it's the OS's responsibility to compile the program, load it into memory and run it.

D. OS needs to be aware of possible errors and take appropriate actions to ensure correct and consistent computing.

**Answer: C**

**The compiler is an application program which does not belong to OS.**

MC5. During the execution of a system call, the following four events occur:

1. Returning to user mode.

2. Execution of the trap instruction.

3. Passing parameters to the system call.

4. Execution of the system call.

The correct order of execution is \_\_\_\_.

A. 2→3→1→4

B. 2→4→3→1

C. 3→2→4→1

D. 3→4→2→1

**Answer: C**

**See 3.23.**

MC6. Which of the following statement is true about system calls?

A. The APIs are alternate versions of system calls as they can completely replace the system calls in the OS implementation.

B. Programs implemented with APIs can often achieve better performance compared to those implemented directly with system calls.

C. The number of parameters passed to system calls is limited if the stack method is used in parameter passing.

D. In the shared-memory interprocess communication model, processes use system calls to create and gain access to regions of memory owned by other processes.

Answer: D

A. API invokes system calls rather than replacing them. See 2.15.

B. APIs typically provide a higher-level interface to system calls, which can introduce overhead in terms of additional function calls and data marshaling. Therefore, by introducing the API layer, a program may potentially lose certain performance.

C. Stack method does not limit the number of parameters being passed. See 2.20.

MC7. Which of the following statement is true about the system structure?

A. Process scheduling function is often provided by the microkernel in a microkernel operating system.

B. The bi-directional dependency relationship between different layers allows more flexibility in a pure layered operating system.

C. The operating system Darwin is designed in the layered approach.

D. The performance of microkernels can greatly suffer due to the overhead of traversing through multiple layers to obtain an OS service.

Answer: A

A. See 2.48.

B. In the layered approach, the relationship between layers is uni-directional. A layer can only invoke operations on lower-level layers. See 2.45.

C. Darwin is a microkernel operating system. See 2.47.

D. The statement is about the drawback of the layered approach. See 2.46 & 2.49.

## [Processes]

MC8. Which of the following statements is true regarding the relationship between a process and a thread?

A. A thread is a lightweight version of a process, and each thread has its own thread control block.

B. A process is a collection of threads, each with its own address space.

C. A process can only have one thread, while a thread can be associated with multiple processes.

D. A process and a thread are identical concepts with different names used by different operating systems.

Answer: A

A. See 4.15.

B. Threads in the same process share the same address space. See 4.7.

C. A thread can be associated with only one process.

D. Thread and process are different concepts.

MC9. Which of the following statements about the pipe is true?

A. An ordinary pipe supports two-way interprocess communication.

B. The ordinary pipeline can only be accessed by one process at a time.

C. In most Unix systems, a named pipe will be automatically deleted by default when the last process using this pipe closes it.

D. A named pipe is a type of file in the file system in UNIX systems that multiple processes can access.

Answer: D

A. Ordinary pipes only support one-way communication. See 3.58.

B. The father process and child process can access the pipe at the same time. See 3.5

C. In most cases, named pipes continue to exist until they are explicitly deleted. See 3.62.

D. See 3.62.

MC10. Which of the following is true about a C program's address space?

A. The program uses the stack to allocate and deallocate memory dynamically.

B. The text section contains the program's executable code.

C. The program counter stores the return address of a function call.

D. The data section is read-only and cannot be modified at runtime.

Answer: B

A. The heap contains memory dynamically allocated. See 3.5.

B. See 3.5.

C. The program counter contains the location of the instruction to be fetched and next executed. See 3.10.

D. The value of global variables stored in the data section can be modified at runtime.

MC11. Which of the following lines is not included in the output of the following C code fragment?

(suppose all *fork()* are successful and necessary header files are included)

```
int a = -1, b = -1;
a = fork();
if (a == 0) b = fork();
a = a > 0 ? 1 : a;
b = b > 0 ? 1 : b;
printf("%d, %d\n", a, b);
fflush(stdout);
```

A. 1, 1

B. 1, -1

C. 0, 1

D. 0, 0

Answer: A

MC12. How many processes will be running concurrently after the *while()* loop?

(suppose only one process runs before the loop, all *fork()* are successful, and necessary header files are included.). Assume *n* is a non-negative integer.

```
int i = n;
while (i--){
    if (fork()) fork();
}
```

- A.  $2^n$
- B.  $3^n$
- C.  $4^n$
- D.  $2n$

Answer: B

### [Threads]

MC13. Which of the following statements about *clone()* and *fork()* in Linux is false?

- A. *fork()* creates an exact copy of the calling process, while *clone()* allows for more control over the creation of the new process or thread.
- B. With the flag `CLONE_FILES`, *clone()* creates a child task that shares the file descriptors with the parent.
- C. *fork()* only works in single-threaded programs, while *clone()* can be used in multi-threaded programs.
- D. All of the above statements are correct.

Answer: C

*fork()* can also be used in multi-threaded programs, and the child process will inherit the information from the calling thread. However, considering the unexpected results, it is not recommended in practical programming.

MC14. Which of the following statements about user threads and kernel threads is true?

- A. User threads are managed and scheduled with kernel support.
- B. The operating system creates a thread control block for each user thread.
- C. Each user thread maps to a kernel thread in the two-level multithreading model.
- D. Blocking of one user thread does not influence the execution of other user threads.

Answer: C

- A. User threads are managed entirely by the user-level threads library. See 4.19 and 4.20.
- B. The operating system only creates TCB for each of the kernel threads. However, there is not always a one-to-one mapping relationship between the user and kernel threads.
- C. See 4.24.
- D. In the many-to-one model, one thread blocking causes all threads mapped to this thread to block. See 4.21.

MC15. Suppose an application is 60% parallel and 40% serial. According to Amdahl's Law, what is the theoretical minimum number of processors required to achieve a speedup of 1.95 times compared to single-core implementation?

- A. 7
- B. 6
- C. 5
- D. 4

Answer: B

Speedup with 5 cores =  $1 / (0.4 + 0.6/5) = 1.923 < 1.95$

Speedup with 6 cores =  $1 / (0.4 + 0.6/6) = 2 > 1.95$

Therefore, we need at least six cores to achieve the speedup of 1.95.

### [CPU Scheduling]

MC16. The convoy effect is associated with \_\_\_\_\_.

- A. FCFS
- B. SJF
- C. Priority
- D. Multi-level queue

Answer: A

MC17. Which of the following scheduling algorithms gives minimum average waiting time?

- A. FCFS
- B. Round Robin
- C. Multi-level feedback queue
- D. SJF

Answer: D

MC18. In multilevel feedback queue algorithm, which of the following statement is true?

- A. Prior knowledge on the next CPU burst time is required for scheduling.
- B. It handles interactive jobs well by delivering similar performance as SJF.
- C. Classification of ready queue is permanent.
- D. It won't suffer from starvation problem.

Answer: B

MC19. Which of the following is true of the rate-monotonic (RM) and earliest-deadline-first (EDF) scheduling algorithms?

- A. RM uses a dynamic priority policy.
- B. Both RM and EDF are non-preemptive.
- C. The rationale for RM is to assign a higher priority to tasks requiring CPU more often.
- D. EDF uses a static priority policy.

Answer: C

MC20. We have 6 processes arriving at the ready queue at the same time. Their estimated runtimes are: 1, 2, 3, 4, 5 and 6 seconds. Their priorities are 3, 1, 4, 2, 5, and 4. A smaller number means higher priority (e.g., priority 1 has a higher priority than priority 2). Ignoring process switching overhead, calculate the average waiting time for priority scheduling.

- A. 5.5s
- B. 8.67s
- C. 8s
- D. 7.67s

Answer: Cancelled. We made some mistakes in the calculation steps.  
The correct steps are as follows:

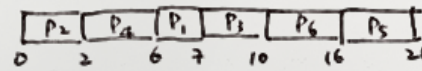
The steps to calculate :

	<u>Burst Time</u>	<u>Priority</u>
P1	1	3 ✓
P2	2	1 ✓
P3	3	4 ✓
P4	4	2 ✓
P5	5	5 ✓
P6	6	4 ✓

Note: P3 and P6 are tied (with priority = 4)

P3 should execute first because it has a smaller ID.

Gantt Chart:



Waiting time:

$$P_1 = 6$$

$$P_2 = 0$$

$$P_3 = 7$$

$$P_4 = 2$$

$$P_5 = 16$$

$$P_6 = 10$$

Average waiting time =

$$\frac{6 + 0 + 7 + 2 + 16 + 10}{6} = \frac{41}{6} = 6.83$$

### [Synchronization]

MC21. Which of the following statements about race condition is true?

- A. A race condition results when several threads try to access and modify the same data separately
- B. A race condition results when several threads try to access the same data concurrently
- C. A race condition will result only if the outcome of execution does not depend on the order in which instructions are executed
- D. None of the above

Answer: D

MC22. Which of the following statements is false?

- A. Spinlocks can be used to prevent busy waits in semaphore execution.
- B. Two standard operations to modify Semaphore: wait() and signal().
- C. Counting semaphores can be used to control access to resources with a limited number of instances.
- D. If a semaphore value is negative, its magnitude is the number of processes currently waiting on the semaphore.

Answer: A



MC23: Which of the following is true about semaphore implementation with no busy waiting?\_\_\_\_\_.

- A. When a process executes wait() and finds the semaphore value is positive, it will suspend itself with block().
- B. A suspended process waiting on the semaphore, may be restarted with wakeup() when some other process executes a signal() operation.
- C. With wakeup(), it will remove one of processes in the waiting queue and place it in the waiting queue.
- D. With block(), it will place the process invoking the operation on the ready queue.

**Answer: B**

MC24. Assume a semaphore S, initialized to 1, shared by several processes. Each process must execute wait(S) before entering the critical section and signal(S) afterward. Suppose a process executes in the following manner.

```
signal(S);  
.....  
/* critical section */  
.....  
wait(S);
```

Which of the following is true?

- A. the critical section is well-protected.
- B. it will get stuck, i.e., a deadlock will occur.
- C. more than one processes may be executing in their critical section.
- D. no process can enter critical section.

**Answer: C**

MC25: Three processes use 4 counting semaphores, initialized as S0 = 1, S1 = 0, S2 = 0, S3 = 0.

```
/*Process P0*/  
while(true) {  
    wait(S0);  
    print '1';  
    release(S1);  
    release(S2);  
    release(S3);  
}
```

```
/*Process P1*/  
wait(S1);  
release(S0);
```

```
/*Process P2*/  
wait(S2);  
release(S0);  
wait(S3);  
release(S0);
```

How many times will '1' be printed out?

- A. At least twice
- B. Exactly twice
- C. Exactly three times
- D. Exactly four times

Answer: D.

## Part II. Calculations [75 points]

### 1. [25 Points] Process and Thread

Suppose all *fork()* are successful and necessary header files are included.

1) (6 points) Consider the following program.

```
void *child(void *arg) {
    int *counter = (int *) arg;
    (*counter)++;
    return NULL;
}
int main(int argc, char *argv[]) {
    int counter = 0;
    pthread_t thread;
    pthread_create(&thread, NULL, child, (void *)&counter);
    pthread_join(thread, NULL);
    printf("%d\n", counter);
    fflush(stdout);
    counter = 0;
    pid_t pid = fork();
    if(pid) wait(NULL);
    else{
        child((void *)&counter);
        exit(0);
    }
    printf("%d\n", counter);
    fflush(stdout);
    return 0;
}
```

What is the output of the program (4 points)? Briefly explain your answer (2 points).

Output:

1  
0

(4 points, 2 points for each line)

Explanation:

The `pthread_create()` will increment `counter` by 1 in both the child thread and the parent thread because the program passes the pointer of the variable `counter` to the child thread. Therefore, the changes of `counter` in the child thread will be reserved in the main thread (1 point).

After the `fork()` call, the data change in the child process only occurs in its own memory because different processes have separate address spaces. Thus, the `counter` in the parent process remains 0 (1 point).

2) (9 points) Consider the following program.

```
void recursion(pid_t pid, int counter){
    if (counter == 0) return;
    if (pid) recursion(fork(), counter - 1);
    printf("%d\n", counter);
    fflush(stdout);
}

int main() {
    int n;
    scanf("%d", &n);
    recursion(fork(), n);
    return 0;
}
```

(a) Given an integer  $n$  ( $n > 0$ ), what is the total number of processes (including the parent process and all the children processes created during execution) (2 points)? Please explain (1 point).

You should write your answer in terms of  $n$ .

**Answer:  $n+2$  (2 points)**

**Explanation:** Every time function `recursion()` is called, a new child process is created. Since the function is called once per iteration and there are  $n+1$  iterations in total, there are  $n+1$  children processes. Therefore, there are  $n+2$  processes in total (1 point).

For the answer, 2 points for  $n+2$  (i.e., the correct answer), 1 point for  $n+1$  (one client process miscounted), 0 point for the others.

For the explanation, 1 point if the explanation is reasonable, otherwise 0 point.

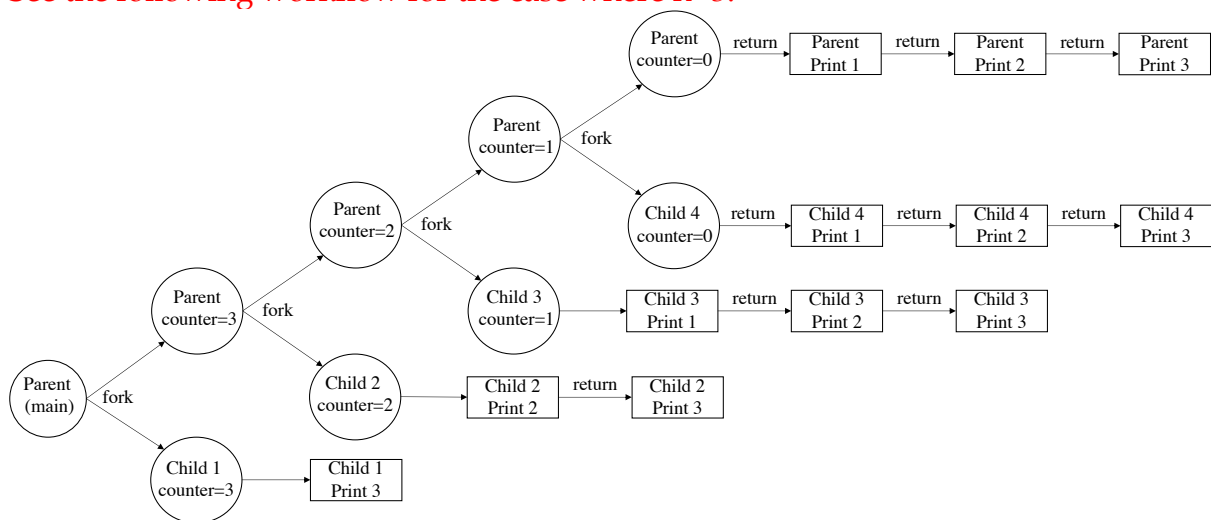
(b) (6 points) How many times are each of the numbers “1”, “2”, and “3” printed when the input number  $n=1, 2$ , and  $3$ , respectively? Write your answer in the following table.

Input value of the variable $n$	Number of “1”s in the output	Number of “2”s in the output	Number of “3”s in the output
1	3	0	0
2	3	4	0
3	3	4	5

1 point for each blank.

Explanation:

See the following workflow for the case where  $n=3$ .



3) (10 points) First, we introduce the function `waitpid()`.

```
pid_t waitpid(pid_t pid, int *status_ptr, int options);
```

This function is similar to the `wait()` function. When the parameter `options=0`, it waits for the termination of the specific child process whose process ID is equal to the parameter `pid`.

In this question, the input comes from three files ("1.txt", "2.txt" and "3.txt") which exist in the same directory as the executable file. Each file includes two lines of input. The first line is a number indicating the length of the vector in the file. The second line is the vector, including multiple integers (See examples for details).

We intend to employ `fork()`, `pipe()`, and input redirection to compute the sum of all numbers in the vectors across the three files with multiprocessing.

Please fill in the following blanks to complete the following program.

**Requirements:**

1. You cannot use system calls except for `close` and `dup` (e.g., `dup2` can't be used).
2. You can write at most 20 characters in each blank.

Here are examples of reading/writing an integer from a file:

```
read(file_id, address_of_int_variable, sizeof(int));
write(file_id, address_of_int_variable, sizeof(int));
```

```
#define SIZE 3
char file_name[SIZE][10] = {"1.txt", "2.txt", "3.txt"};
int main() {
    int fd[SIZE - 1][2], i, j;
    pid_t pid[SIZE - 1];
    for(i = 0; i < SIZE - 1; i++) {
        pipe(fd[i]);
        pid[i] = fork(); // Record the pid of the child
        if(!pid[i]) break;
    }
    int fin = open(file_name[i], O_RDONLY, S_IRUSR | S_IWUSR);
    int length, number, sum = 0;
    _____ Blank 1 _____; // Close the standard input
    _____ Blank 2 _____; // Input redirecting
    scanf("%d", &length); // Input length of the vector
    for (j = 0; j < length; j++){
        scanf("%d", &number); // Input the vector
        _____ Blank 3 _____;
    }
    if(!pid[i]) {
        close(fd[i][0]); // Close read end
        write(_____ Blank 4.1 _____, _____ Blank 4.2 _____, sizeof(int)); // Send data
        exit(0);
    }
    for(i = 0; i < SIZE - 1; i++) {
        int status, child_sum;
        waitpid(pid[i], &status, 0); // Wait for child process
```

```

        close(fd[i][1]); // Close write end
        read(____ Blank 5.1 ____, ____ Blank5.2 ____, sizeof(int)); // Receive data
        ____ Blank 6 ____; // summing up
    }
    printf("Summation is %d.\n", sum);
    return 0;
}

```

Sample Input and Output.

1.txt	2.txt	3.txt	Output in the console
4 20 10 30 20	2 40 30	3 20 10 50	Summation is 230.

Explanation of the sample:  $230 = 20+10+30+20+40+30+20+10+50$ .

Write your answer in the following blanks.

**Answer:**

Blank 1: close(0) or close(fileno(stdin)) or close(STDIN\_FILENO) (1 point)

Blank 2: dup(fin) (1 point)

Blank 3: sum += number (2 points)

Blank 4.1: fd[i][1] (1 point)

Blank4.2: &sum (1 point)

Blank 5.1: fd[i][0] (1 point)

Blank5.2: &child sum (1 point)

Blank 6: sum += child sum (2 points)

For blank 5.2 and 6, “child\_sum” can be replaced by some other variables like “number”. But the students should use the same variable in these two blanks.

## 2. [35 Points] CPU Scheduling

1) (13 points) The following table contains 5 processes (P1-P5). Each process has its arrival time, burst time, and priority.

Process	Arrival Time	Burst Time	Priority (Only used in Priority Scheduling)
P <sub>1</sub>	0	5	2
P <sub>2</sub>	1	3	1
P <sub>3</sub>	4	8	2
P <sub>4</sub>	8	2	1
P <sub>5</sub>	10	6	2

- Whenever there is a tie among processes (the same arrival time, remaining time, priority, etc), they are inserted into the ready queue in the ascending order of process id.
- A smaller number means a higher priority (e.g., priority 1 has a higher priority than priority 2)

For each of the following scheduling algorithms, draw the Gantt charts depicting the sequence of the process execution and calculate the average turnaround time.

1. FCFS
2. SJF (non-preemptive)
3. Preemptive priority scheduling with RR (with quantum of 2 milliseconds)
4. SRTF



- **FCFS (3 points)**

P1	P2	P3	P4	P5	
0	5	8	16	18	24

Avg turnaround:  $(5 + 7 + 12 + 10 + 14)/5 = 9.6$

Partial credits:

- The Gantt chart is correct up to  $t=8$  (1 mark)
- The Gantt chart is correct (2 marks)
- The Gantt chart and the average turnaround time are both correct (3 marks)

- **SJF (3 points)**

P1	P2	P4	P5	P3	
0	5	8	10	16	24

Avg turnaround:  $(5 + 7 + 2 + 6 + 20)/5 = 8$

Partial credits:

- The Gantt chart is correct up to  $t=8$  (1 mark)
- The Gantt chart is correct (2 marks)
- The Gantt chart and the average turnaround time are both correct (3 marks)

- **Preemptive priority scheduling with RR (TQ=2) (4 points)**

P1	P2	P1	P3	P1	P4	P1	P3	P5	P1	P3	P5	P3	P5	
0	1	4	5	7	8	10	11	13	15	16	18	20	22	24

Avg turnaround:  $(16 + 3 + 18 + 2 + 14) / 5 = 10.6$

Partial credits:

- The Gantt chart is correct up to  $t=8$  (1 mark)
- The Gantt chart is correct up to  $t=16$  (2 marks)
- The Gantt chart is correct (3 marks)
- The Gantt chart and the average turnaround time are both correct (4 marks)

- **SRJF (3 points)**

P1	P2	P1	P4	P5	P3	
0	1	4	8	10	16	24

Avg turnaround:  $(8 + 3 + 20 + 2 + 6) / 5 = 7.8$

Partial credits:

- The Gantt chart is correct up to  $t=8$  (1 mark)
- The Gantt chart is correct (2 marks)
- The Gantt chart and the average turnaround time are both correct (3 marks)

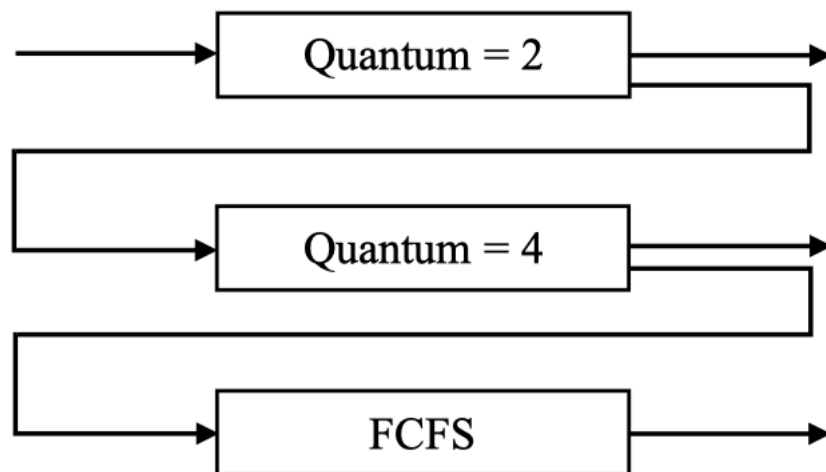
2) (10 points) Consider the following single-thread processes, arrival times, burst times and the following three queues:

Q0 - RR with time quantum 2 milliseconds

Q1 - RR with time quantum 4 milliseconds

Q2 - FCFS

Process	Arrival Time	Burst Time
P <sub>1</sub>	0	5
P <sub>2</sub>	4	13
P <sub>3</sub>	6	7
P <sub>4</sub>	11	4
P <sub>5</sub>	14	14
P <sub>6</sub>	20	6



a) (6 points) Draw the Gantt chart depicting the scheduling procedures for these processes.

P1	P1	P2	P3	P1	P2	P4	P2	P5	P2	P3	P6	P3	P4	P5	P6	P2	P3	P5	
0	2	4	6	8	9	11	13	14	16	17	20	22	23	25	29	33	40	41	49

Partial credits:

- Too many mistakes (0 marks)
- Understanding the basics with several mistakes (2 marks)
- Almost correct, with one minor mistake (4 marks)
- All correct (6 marks)

b) (4 points) Calculate the average waiting time.

P1: 4

P2: 23

P3: 28

P4: 10

P5: 21

P6: 7

average waiting time =  $(4+23+28+10+21+7)/6=15.5$

Partial credits:

- 1-3 marks for the steps (depending on the errors in the steps)
- 4 marks for the final answer is correct (or the answer and steps are correct)

- 3) (12 points) Consider the following single-thread processes, executing times, deadlines and periods. Assume all processes arrive at timeslot 0. Fill in the table with the ID of the process that is running on the CPU with Rate-Monotonic (RM) scheduling and Earliest Deadline First (EDF) scheduling in the first 20 timeslots, and show how many deadlines are missed in each scheduler.

Process	Processing Time	Deadline	Period
P <sub>1</sub>	1	4	5
P <sub>2</sub>	2	5	10
P <sub>3</sub>	4	13	16
P <sub>4</sub>	2	7	8

RM																					
EDF																					
Time	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

How many deadlines are missed for RM: \_\_\_\_\_

How many deadlines are missed for EDF: \_\_\_\_\_

**Solution:**

RM	1	4	4	2	2	1	3	3	4	4	1	2	2	3	3	1	4	4	3	3
EDF	1	2	2	4	4	1	3	3	3	3	1	2	2	4	4	1	4	4	3	3
Another EDF Solution	1	2	2	4	4	1	3	3	3	3	1	4	4	2	2	1	4	4	3	3
Time 0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

RM: 1(P<sub>3</sub>), EDF: 0 (2 marks)

RM: (4 marks), EDF: (6 marks)

Partial Marking scheme for Q2.3:

- 4 marks for timeslots of RM, 2 marks for just one minor mistake, 0 otherwise.
- 4 marks for timeslots of EDF, 2 marks for just one minor mistake, 0 otherwise.
- 4 marks for missing deadline questions, 2 for each one.

### 3. [15 points] Synchronization

- 1) (5 points) Suppose two threads execute the following C code concurrently and access shared variables `a`, `b`, and `c`.

Initialization:

```
int a = 1;
int b = 3;
int c = 0;
```

// Thread 1	// Thread 2
a = b + 1	a = 13
c = a - b	b = 10
	c = 2

What are all the possible values for `c` after both threads complete? You can assume that readings and writings of the variables are atomic, and that the order of execution of statements within each thread is preserved by the C compiler and the hardware, so it matches the code above.

Answer: -6, 1, 2, 3, 10 (5 points)

1 point for each correct answer.

1 point deduction for each incorrect output (until 0 point for this question).

Explanation:

-6:

Thread 1 executes  $a=3+1=4$ . After thread 1 loads  $a=4$  into the register for the subtraction, thread 2 executes all its statements which leads to  $b=10$ . Then thread 1 loads  $b=10$  into the register. At last, thread 1 calculates  $c=4-10=-6$ .

1:

Case 1: Thread 1 executes  $a=b+1=4$ . After thread 1 loads 4 and 3 into the registers for the subtraction, thread 2 executes all its statements. And then thread 1 runs  $c=4-3=1$ .

Case 2: Thread 2 executes  $a=13$ ,  $b=10$ , and  $c=2$ . Then, thread 1 runs  $a=10+1=11$  and calculates  $c=11-10=1$ .

2:

The last statement of thread 2 " $c=2$ " is executed at last. Therefore,  $c=2$ .

3:

Thread 1 executes  $a=3+1=4$ . Then thread 2 executes  $a=13$ ,  $b=10$ , and  $c=2$ . After that, thread 1 executes  $c=13-10=3$ .

10:

**A lot of students missed this value.** Thread 1 executes  $a=3+1=4$ . Then thread 2 executes  $a=13$  which overwrites the change by thread 1. After that, thread 1 **loads  $a=13$**  and  **$b=3$**  for the second statement. Before thread 1 finishes the subtraction, thread 2

executes all its remaining statements. Thread 1 calculates  $c=13-3=10$  and stores the value in  $c$ .

Note:

A lot of students write  $c = -9$  which is not possible.  $-9$  comes from  $a = 1$  and  $b = 10$ . But before this,  $a$  is already modified by either  $a = b+1$  or  $a = 13$  in the two threads. So this case is not possible.

2) (10 points) Barrier: A common parallel programming pattern is to perform processing in a sequence of parallel stages: all threads work independently during each stage, but they must synchronize at the end of each stage at a synchronization point called a **barrier**. If a thread reaches the barrier before all other threads have arrived, it waits. When all threads reach the barrier, they are notified and can begin the execution on the next phase of the computation.

Example:

```
while (true) {
    Compute stuff;
    BARRIER();
    Read other threads results;
}
```

- a) (4 points) The following implementation of Barrier is incomplete and has four lines missing. Fill in the missing lines so that the Barrier works according to the prior specifications.

Hint: You can use the following function calls to fill in the missing blanks:

```
CV.signal()
CV.signalAll()
CV.wait()
L.acquire()
L.release()
```

```
class Barrier {
private:
    int numWaiting {0}; // Initially, no one at barrier
    int numExpected {0}; // Initially, no one expected
    Lock L;
    ConditionVar CV;

public:
    void threadCreated() {
        L.acquire();
        numExpected++;
        L.release();
    }
    void enterBarrier() {
        _____ Blank1 _____;
        numWaiting++;
        if (numExpected == numWaiting) { // If we are the last
            numWaiting = 0; // Reset barrier and wake threads
            _____ Blank2 _____;
        } else { // Else, put me to sleep
            _____ Blank3 _____;
        }
        _____ Blank4 _____;
    }
};
```

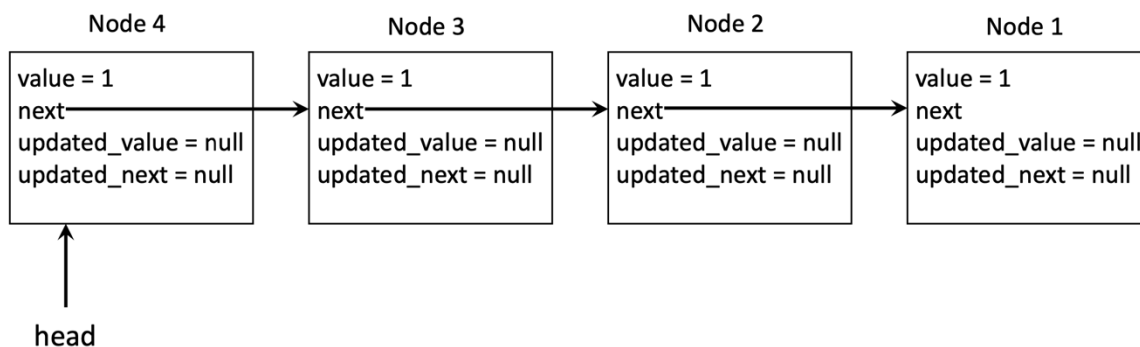


Blank1	
Blank2	
Blank3	
Blank4	

Solution:

Blank1	<u><i>L.acquire()</i></u>
Blank2	<u><i>CV.signalAll()</i></u>
Blank3	<u><i>CV.wait()</i></u>
Blank4	<u><i>L.release()</i></u>

- b) (6 points) Now, let us use the above Barrier implementation in a parallel algorithm. Consider the linked list below:



In our parallel algorithm, there are four threads (Thread 1, Thread 2, Thread 3, Thread 4). Each thread has its own instance variable *node*, and all threads share the class variable *barrier*. Initially, Thread 1's *node* references Node 1, Thread 2's *node* references Node 2, Thread 3's *node* references Node 3, and Thread 4's *node* references Node 4.

In the initialization steps, **barrier.threadCreated()** is called once for each thread created, so we have **barrier.numExpected == 4** as a starting condition.

Once all four threads are initialized, each thread calls its **run()** method. The **run()** method is identical for all threads:

```

void run() {
    bool should_print = true;
    while (true) {
        if (node.next != NULL) {
            node.updated_value = node.value + node.next->value;
            node.updated_next = node.next->next;
        } else if (should_print) {
            std::cout << node.value;
            should_print = false;
        }
        barrier.enterBarrier();
        node.value = node.updated_value;
        node.next = node.updated_next;
        barrier.enterBarrier();
    }
}
  
```

List all the values that are printed to stdout along with the thread that prints each value. For example, "Thread 1 prints xxx".

Answer:

Thread 1 prints 1 (1 points)

Thread 2 prints 2 (1 points)

Thread 3 prints 3 or "Thread 3 prints 2 + null" or any other answer with an explicit explanation of what happens when you add a number with null (2 points). "Threads 3 prints 2" only gets 1 point.

Thread 4 prints 4 (2 points)

Note: This is a parallel list ranking algorithm where the final value of each node is its position in the linked list.