

University of California, Berkeley  
 College of Engineering  
 Computer Science Division – EECS

Spring 2017

Ion Stoica

### Third Midterm Exam

April 24, 2017

CS162 Operating Systems

<b>Your Name:</b>	
<b>SID AND 162 Login:</b>	
<b>TA Name:</b>	
<b>Discussion Section Time:</b>	

General Information:

This is a **closed book and one 2-sided handwritten note** examination. You have 80 minutes to answer as many questions as possible. The number in parentheses at the beginning of each question indicates the number of points for that question. You should read **all** of the questions before starting the exam, as some of the questions are substantially more time consuming.

Write all of your answers directly on this paper. ***Make your answers as concise as possible.*** If there is something in a question that you believe is open to interpretation, then please ask us about it!

**Good Luck!!**

QUESTION	POINTS ASSIGNED	POINTS OBTAINED
<b>1</b>	<b>24</b>	
<b>2</b>	<b>15</b>	
<b>3</b>	<b>12</b>	
<b>4</b>	<b>16</b>	
<b>5</b>	<b>15</b>	
<b>6</b>	<b>18</b>	
<b>TOTAL</b>	<b>100</b>	

**P1** (24 points total) True/False and Why? **CIRCLE YOUR ANSWER.** For each question: 1 point for true/false correct, 2 point for explanation. An explanation cannot exceed 2 sentences.

- a) When two processes `mmap()` the same region of a file, `mmap()` will always return the same set of virtual addresses to both processes, which represent a shared region of memory.

TRUE

FALSE

Why?

The virtual addresses are not necessarily the same, both sets will point to the same physical region.

- b) A TCP sender can overflow the receiver's buffer.

TRUE

FALSE

Why?

Flow control will make sure the sender never sends more bytes than the receiver can store.

- c) It is not possible to develop an agreement protocol between two hosts (i.e., a protocol in which the two hosts agree to do something simultaneously) over an unreliable network?

**TRUE**

**FALSE**

Why?

Assume there exists a protocol, and assume  $N$  is the minimum number of message of any such protocol. Since the sender of the last message cannot be sure whether the message has been delivered, we can remove the message from the protocol, which contradicts the fact that  $N$  is the smallest number of messages. This is the General's Paradox.

- d) Two-Phase Commit does not guarantee all ACID properties (Atomicity, Consistency, Isolation, and Durability).

**TRUE**

**FALSE**

Why?

Only has Atomicity and Durability.

- e) The end-to-end principle states that one should make the lower levels between the higher-level ends as reliable and secure as possible.

**TRUE**

**FALSE**

Why?

The E2E principle says the opposite.

- f) FAT provides slower sequential access to a file than Unix File System.

Why? **TRUE**

FALSE

FAT needs to perform a disk access for every block to find the pointer to that block, while Unix File System gets the pointers to many blocks by reading an inode.

- g) Using quorum consensus a writer can return before the data has been written to all replicas.

Why? **TRUE**

FALSE

The writer returns after it hears from  $W$  replicas, where  $W$  is typically smaller than the total number of replicas,  $N$ .

- h) With consistent hashing (e.g., Chord) if a new node is added to a system consisting of  $N > 1$  nodes, this will lead to keys from every node to be moved around.

Why? TRUE

**FALSE**

Let's  $A$  be the new node. Only the keys of  $A$ 's successor will be moved to  $A$ .

**P2 (15 points) Transactions.** Consider two variables  $A$  and  $B$  which are initialized to 1, and 2, respectively:

A	1
B	2

Consider the following two transactions being attempted on A and B:

Transaction 1			Transaction 2		
A=3	B=4	Commit	A=5	B=6	Commit

For the following scenarios, list which ACID property was violated or “none” if the transactions obeyed ACID, explain your answer in 1 sentence:

a) (3 points) Transaction 1 successfully committed, then the computer crashed after step 2 (B=6) of transaction 2 (but before writing the commit record to the log).

State after Recovery:

A	3
B	4

Correct ACID behavior. Txn 1 committed, txn 2 did not.

b) (3 points) The transactions were run on two threads (at the same time) and have both committed. Note: concurrent transactions may complete in any order. State which order your answer assumes (if any).

State after Recovery:

A	3
B	4

Order (2, 1): Correct ACID behavior. Isolation only requires some serialized order (the outcome is as if one transaction ran to completion, then the next ran to completion), overwriting old commits is perfectly acceptable.

Order (1, 2): Durability (or maybe isolation) was violated. Durability: The changes from the latest committed transaction were not reflected in the final state (most recent changes were lost). Isolation is possible if you assume that B logged its changes, then A ran to completion, then B committed.

c) (3 points) The transactions were run on two threads (at the same time) and have both committed. Note: concurrent transactions may complete in any order. State which order your answer assumes (if any).

Final State:

A	3
B	6

Isolation was violated. A and B got interleaved. In this case, there was no total “order” so they may assume either order.

d) (3 points) **Transaction 1 successfully committed**, then the computer crashed during transaction 2 (before commit).

State after Recovery:

A	5
B	4

Atomicity was violated. B was only partially performed.

e) (3 points) Both transactions committed (in order txn1,txn2), then the computer crashed.

State after Recovery:

A	3
B	4

Durability was violated. Txn 2 was lost despite committing.

**P3 (12 points) Synchronization.** Consider the following C program.

```
int thread_count = 0;

void *thread_start(void *arg) {
    thread_count++;
    if (thread_count == 3) {
        char *argv[] = {"/bin/echo", "162 is the best!", NULL};
        /* Replace current process image with a new process image
         * that runs program pointed by "*argv" with arguments "argv" */
        execv(*argv, argv);
    }
    printf("Thread: %d\n", thread_count);
    return NULL;
}

int main(int argc, char *argv[]) {
    int i, status;
    for (i = 0; i < 2; i++) {
        pid_t pid = fork();
        if (pid != 0) {
            thread_count++;
            wait(&status);
        }
        pthread_t *thread = malloc(sizeof(pthread_t));
        /* start new "thread". The new thread starts execution by invoking
         * "thread_start()" with NULL arg. */
        pthread_create(thread, NULL, &thread_start, NULL);
        pthread_join(*thread, NULL);
    }
    return 0;
}
```

a) (8 points) Provide one possible output from this C program.

```
Thread: 1
Thread: 2
162 is the best!
Thread: 2
162 is the best!
Thread: 4
```

b) (4 points) Is the output deterministic? If yes, give a short explanation why. If no, explain how to make the output deterministic.

Yes. Each process is serialized by its wait, and each thread is serialized by `pthread_join`.

**P4 (16 points) Queueing theory** Consider a system with a single queue and a single server. The arrival rate of the requests follows a Poisson process, and the service time is exponential distributed. Thus, the system implements an M/M/1 queue. Assume that:

- The utilization  $U$  is 50%.
- The average service time  $T_{ser} = 1/6$  seconds.

a) (4 points) What is the average service rate  $\mu$  and arrival rate  $\lambda$ ?

Service rate  $\mu$  is  $1/T_{ser}$ , so 6 jobs per second.

Utilization,  $u = \lambda/\mu$ , so  $\lambda = u * \mu = 3$  jobs per second.

$\lambda$

a) (4 points) What are the average time spent in the queue  $T_q$  and the response time?

$T_q = T_{ser} * u / (1-u) = (1/6) * (1/2) / (1-1/2) = 1/6$  seconds.

$T_{sys} = T_q + T_{ser} = 1/6 + 1/6 = 1/3$  seconds.

b) (4 points) What is the average length of the queue  $L_q$ ?

By little's law, this is  $\lambda * T_q = 3 * 1/6 = 1/2$  jobs.

c) (4 points) Assume that we add an additional server with the same service time ( $T_{ser} = 1/6$  seconds) to help service our queue. What is the average time spent in the queue  $T_q$  and response time for the new system, assuming the same arrival rate as at point (a), and that requests are balanced equally across the two servers?

$u$  is now half because the requests are split over the two servers.

$T_q = T_{ser} * (u/(1-u)) = (1/6) * (1/4) / (1-1/4) = 1/18$  seconds.

$T_{sys} = T_q + T_{ser} = 1/18 + 1/6 = 4/18$  seconds.



**P5: (15 points) Caching and storage:** You have been hired at a hot startup located in downtown Berkeley and you are tasked with building a key-value store backed by spinning disks. Each disk has the following performance:

- Bandwidth: 100MB/s
- Disk controller latency: 1ms
- Average rotational latency: 3ms
- Seek latency: 3ms

In front of the backing store you have provisioned an in-memory caching layer. The performance of this layer is bounded by the capacity of the network, which is 100MB/s.

Each access to the key-value requests translates into one access into the cache and, if there is a cache miss, to one access to one backing disk. Each access results in reading 1 key-value of 100 KB, which is stored contiguously on the disk (i.e., on adjacent sectors.) Ignore the network and cache latency.

- a) (5 points) Considering your workload hits the cache 50% of the time, what is the average service time (i.e., the time/latency to retrieve a key-value)?

Service time to read a key-value from disk =  
 = controller latency + rotation latency + seek latency + transfer time  
 = 1ms + 3ms + 3ms + 100KB/100MBps = 8ms

Service time to read a key-value from cache: 100KB/100MBps = 1ms

Average service time =  $0.5 \cdot 1\text{ms} + 0.5 \cdot (8\text{ms} + 1\text{ms}) = 0.5 + 4.5 = 5\text{ms}$ .

We gave partial credit for not including service time to read a key-value from cache.

We also considered solutions that don't account for cache updates as correct (average 4.5ms).

- b) (5 points) What is the maximum number of requests (reads) per second that this system can handle?

Requests per second =  $1\text{s}/5\text{ms} = 200 \text{ request/sec}$

We gave partial credit to solutions that didn't include service time to read a key-value from cache.

- c) (5 points) Now consider that you decide to change your spinning disks with faster SSDs to achieve an average service time (i.e., the time/latency to retrieve a key-value) of 1.25ms. Assume that the SSD controller latency is 1ms, the same as the HDD controller latency. What is the bandwidth of the SSD?

Service time to read from SSD =  $(1\text{ms} + 100\text{KB}/X)$

Average service time =  $0.5 * 1\text{ms} + 0.5 * (1\text{ms} + 100\text{KB}/X) = 1.25\text{ms} \Rightarrow$   
 $1\text{ms} + 100 \text{KB}/X = 1.5\text{ms} \Rightarrow X = 200\text{MB/s}$

We gave partial credit to solutions that didn't include service time to read a key-value from cache.

**P6 (18 points) TCP Flow Control.**

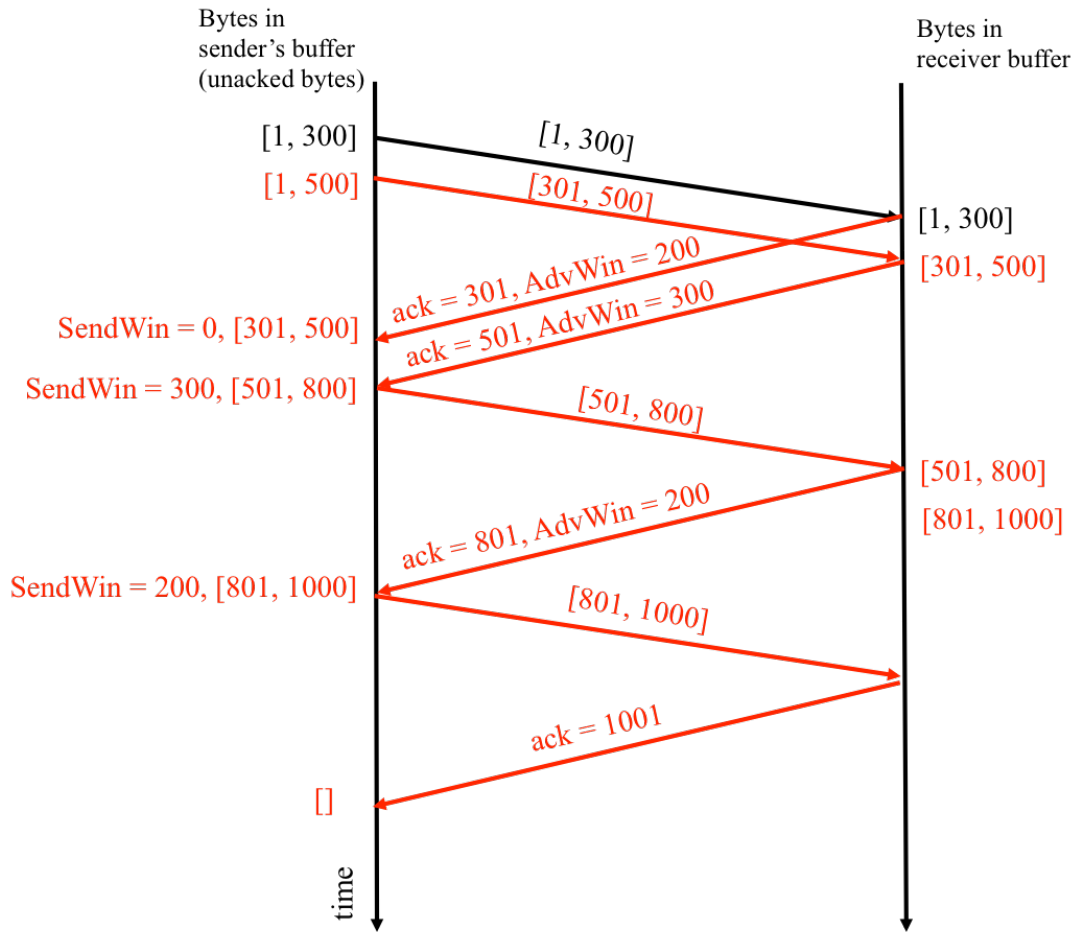
Assume host S uses TCP to send 1000 bytes of data to host R. We make the following assumptions:

- The maximum packet size is 300B.
- Maximum size of R's receiving buffer is 500B.
- The delay between the machines is much larger than the time it takes to S to send 1KB, i.e, S could send 1KB faster than it takes to hear back an ack from R.
- R consumes the bytes as soon as it receives them (i.e., R's OS sends data to receiving app as soon as it gets them.)

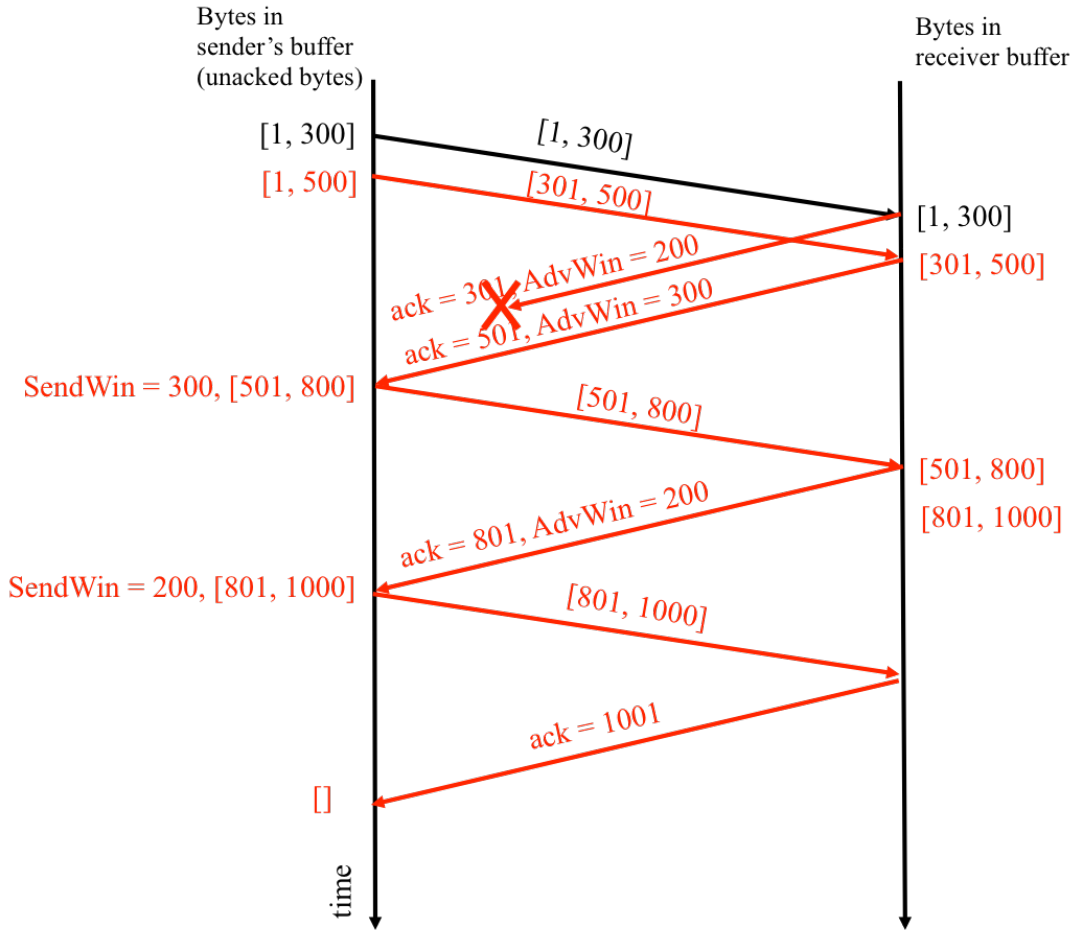
a) (6 points) Assume no packet or acknowledgement is lost. Draw the packet diagram also indicating the bytes in R's sending buffer and the bytes in S's receiving buffer. For the acknowledgements indicate both the sequence of next expected byte, and the advertised window.

The first packet with the associated information is drawn for you.

Initial advertised window (AdvWindow) is equal with the receiver buffer size which is 500B. Also, while below we also show sender window (SendWin), there was **no** penalty for not showing it, as the problem didn't ask for it.



b) (6 points) Draw the packet diagram assuming the acknowledgement of the first packet is lost. No other packet or acknowledgement is lost.



c) (6 points) Draw the packet diagram assuming the third packet from the sender is lost. No other packet or acknowledgement is lost.

