



Department of Electrical Engineering and Computer Science

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

6.824 Distributed System Engineering: Spring 2017

Exam I

Write your name on this cover sheet. If you tear out any sheets, please write your name on them. You have 80 minutes to complete this quiz.

Some questions may be much harder than others. Read them all through first and attack them in the order that allows you to make the most progress. If you find a question ambiguous, write down any assumptions you make. Write neatly. In order to receive full credit you must answer each question as precisely as possible.

You may use class notes, papers, and lab material. You may read them on your laptop, but you are not allowed to use any network. For example, you may not look at web sites or communicate with anyone.

I (7)	II (32)	III (21)	IV (7)	V (22)	VI (8)	VII (3)	Total (100)

Name:

I Lab 1

Alyssa P. Hacker is writing a simplified version of the `schedule` function for the MapReduce lab. This function takes in an array of `Tasks` for a particular map or reduce phase and sends them to available workers. The function should return once all of the tasks in this phase have been completed.

```
01 func schedule(tasks []Task) {
02     var mu sync.Mutex
03     tasksRunning := 0
04     for _, task := range tasks {
05         go func(task Task) {
06             // Count this task as running.
07             mu.Lock()
08             tasksRunning++
09             mu.Unlock()
10
11             // Get an available worker and let them run the task.
12             worker := <-registerChan
13             call(worker, "Worker.DoTask", task)
14
15             // Indicate that this task has completed.
16             mu.Lock()
17             tasksRunning--
18             mu.Unlock()
19
20             // Reregister the worker as available.
21             registerChan <- worker
22         }(task)
23     }
24     // Wait for all the tasks to be done.
25     for {
26         mu.Lock()
27         t := tasksRunning
28         mu.Unlock()
29         if t == 0 {
30             break
31         }
32     }
33     return
34 }
```

1. [7 points]: Alyssa notices that `schedule` sometimes returns before all of the tasks have been finished, even without worker failures. Explain a sequence of events that could cause `schedule` to return too early.

II Short reading questions

2. [8 points]: Chunk servers in GFS (as described by Ghemawat *et al.*) store the version number of a chunk persistently on disk. How would GFS break if chunk servers stored the version number of each chunk only in memory? Briefly explain your answer.

3. [8 points]: Does Spinnaker's timeline consistency (as described in "Using Paxos to build a scalable, consistent and highly available datastore" by Rao et al.) allow a client to read a stale value? Specifically, can a client read an old value for a key after another client has successfully put a new value for that key? If yes, explain how this can happen; if no, explain why it can't happen.

In VM-FT (as described (by Scales *et al.*), when a virtual machine tries to read the time-of-day clock, the hypervisor on the primary machine obtains control and reads the real hardware time-of-day clock. The hypervisor returns the read value to the virtual machine and forwards the read value to the hypervisor on the backup machine. The hypervisor on the backup returns the value read at the primary to its virtual machine when that virtual machine reads the time-of-day clock.

4. [8 points]: Ben observes that his backup machine has a hardware time-of-day clock too, and he also sees that the time-of-day clocks on the primary and the backup are perfectly synchronized – if you read them both at the same instant, they always return exactly the same value. Ben proposes a simplified VM-FT design that avoids the communication between the primary and the backup for the time-of-day clock. When the VM on the primary wants to read the time-of-day, the hypervisor returns the value from its local clock, and doesn't forward the value to the backup. When the virtual machine on the backup wants to read the time-of-day clock, the hypervisor on the backup reads its local time-of-day clock and returns that value to its virtual machine. Even though the clocks on Ben's machines are perfectly synchronized, Ben observes that the virtual machine on the primary and the backup diverge over time. Briefly explain why Ben's design doesn't work.

5. [8 points]: Ben wants to add the following scheme for duplicate request detection to ZooKeeper. His proposal is to modify Zookeeper to maintain a duplicate detection table, replicated across the Zookeeper servers. Each client stamps each of its requests with a unique ID. After Zab agrees on a request, and before a server executes it, the server indexes its table with the request's unique ID to see whether it has already executed the request. If there is no entry for that unique ID, the server executes the request and stores the result in its table under the unique ID. If the unique ID is present in the table, the server doesn't execute the request, and if it was the server that received the request from the client it sends the result stored in the table as the response to the client. This design has a serious flaw – what is it? (Briefly explain your answer.)

III Raft

Answer each of these questions with reference to the Raft protocol described in the paper's Figure 2 (without membership change, without log compaction, and without Section 8).

6. [7 points]: Suppose the leader waited for **all** servers (rather than just a majority) to have $matchIndex[i] \geq N$ before setting `commitIndex` to N (at the end of Rules for Servers). What specific valuable property of Raft would this change break?

7. [7 points]: Suppose the leader immediately advanced `commitIndex` to its last log entry without waiting for responses to `AppendEntries` RPCs, and regardless of the values in `matchIndex`. Describe a specific situation in which this change would lead to incorrect execution of the service being replicated by Raft.

8. [7 points]: There are five Raft servers. S1 is the leader in term 17. S1 and S2 become isolated in a network partition. S3, S4, and S5 choose S3 as the leader for term 18. S3, S4, and S5 commit and execute commands in term 18. S1 crashes and restarts, and S1 and S2 run many elections, so that their currentTerms are 50. S3 crashes but does not re-start, and an instant later the network partition heals, so that S1, S2, S4, and S5 can communicate. Could S1 become the next leader? How could that happen, or what prevents it from happening? Your answer should refer to specific rules in the Raft paper, and explain how the rules apply in this specific scenario.

IV Lab 2

Ben's Raft is failing the Lab 2 tests, which say that log entries are appearing on the apply channel out of order. Here is the part of his code that sends on the apply channel. Ben's implementation calls this function every time it updates the commit index. There is a single lock (Go mutex) protecting each Raft instance, and the calling code always holds the lock.

```
01. func (rf *Raft) applyToService() {
02.   for rf.commitIndex > rf.lastApplied {
03.     entry := rf.log[rf.lastApplied + 1]
04.     msg := ApplyMsg {
05.       Index: rf.lastApplied + 1,
06.       Command: entry.Command,
07.     }
08.     rf.lastApplied++
09.     go func() { rf.applyCh <- msg } ()
10.   }
11. }
```

9. [7 points]: Explain Ben's bug.

V Concurrency Control

Suppose you have a storage system that provides serializable transactions, and that can abort transactions.

The database holds three objects, x , y , and z , with these initial values:

```
x = 1
y = 2
z = 0
```

Only two transactions execute, and they start at about the same time. This is what they do:

```
T1:
  x=10
  y=20
```

```
T2:
  temp = x + y
  z = temp
```

10. [8 points]: Which final values are possible (after each transaction has either aborted or finished committing and applying all updates)? Mark each of the following as “yes” or “no”.

- a. **Yes / No** $x=1$ $y=2$ $z=0$
- b. **Yes / No** $x=1$ $y=2$ $z=30$
- c. **Yes / No** $x=1$ $y=20$ $z=21$
- d. **Yes / No** $x=10$ $y=20$ $z=3$
- e. **Yes / No** $x=10$ $y=20$ $z=30$

11. [7 points]: Suppose T2 has just executed its first line, but has not tried to commit. If the transaction system uses pessimistic two-phase locking, can the first line of T2 see $x=10$ and $y=2$? How, or why not?

12. [7 points]: Suppose T2 has just executed its first line, but has not tried to commit. If the transaction system uses optimistic concurrency control (e.g. FaRM), can the first line of T2 see $x=10$ and $y=2$? How, or why not?

VI FaRM

Ben Bitdiddle reads the FaRM paper (“No compromises: distributed transactions with consistency, availability, and performance”). Ben thinks the Figure 4 commit protocol looks too much like two phase commit with pessimistic locking. He decides to get rid of the locks, and use only version numbers. His new commit protocol has no LOCK phase, but instead uses VALIDATE for written objects as well as read objects. Note that FaRM causes a transaction to first read every object that it writes (in order to fetch the version number).

Here is Ben’s modified commit protocol:

Phase 1, VALIDATE: The coordinator uses one-sided reads to re-fetch the version numbers of all objects used in the transaction, and aborts if any version has changed since the object was first read by the transaction. (The same as FaRM’s VALIDATE in Figure 4 / Section 4, but used for written objects too.)

Phase 2, COMMIT BACKUP: The same as FaRM.

Phase 3, COMMIT PRIMARY: The coordinator sends the primary the updated object contents in this message. The primary copies the new object content to the object’s location in memory and increments its version number.

Phase 4, TRUNCATE: The same as FaRM, except that the backups don’t unlock.

- 13. [8 points]:** Assume there are no computer or network failures. Describe a specific situation in which transactions would yield incorrect results with Ben’s modified commit protocol.

VII 6.824

14. [1 points]: Lab 2 is challenging. If you ran into significant difficulties, what were they? What can we do to avoid these difficulties for future generations of 6.824 students?

15. [1 points]: Which papers should we definitely keep for future years?

16. [1 points]: Is there any paper we should delete?

End of Exam I