



Department of Electrical Engineering and Computer Science

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

6.5840 Distributed System Engineering: Spring 2024

Exam I

Please write your name on the bottom of each page. 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, use ChatGPT, or communicate with anyone.

Name:

Gradescope E-Mail Address:

I MapReduce

Have a look at Figure 3(a) in the paper *MapReduce: Simplified Data Processing on Large Clusters* by Dean and Ghemawat. The three graphs on the left show the rate of data movement over time for a MapReduce job that sorts a terabyte of data: the rate at which Maps read their input, the rate at which intermediate data is shuffled, and the rate at which Reduces write their output. For these questions you should assume that only this MapReduce job is using the servers and network, and that there are no failures. Many of the numbers below are derived from looking at the graphs, and are thus approximate; your reading of the graphs may be somewhat different from our's; you should circle the answer that is closest to what you think is correct.

1. [6 points]: Roughly when is the first time at which the sort application's Reduce() function is called? Circle the best answer.

- 0 seconds
- 50 seconds
- 150 seconds
- 300 seconds

2. [7 points]: Roughly how long does it take a single application Reduce function to sort its share of the data (just the sort, not including either the shuffle or the writing of the output)? Circle the best answer.

- 10 seconds
- 75 seconds
- 200 seconds
- 250 seconds
- 650 seconds

Name: _____

3. [6 points]: Why are there two bumps in the Shuffle graph? That is, why does the Shuffle graph go up and then down from time 20 to 200, remain at zero for 100 seconds, and then go up and then down from time 300 to 600? Circle the best answer.

- There are more Map tasks ($M = 15,000$) than there are machines.
- There are more Reduce tasks ($R = 4000$) than there are machines.
- There are more Map tasks than there are Reduce tasks.
- The aggregate network throughput is smaller than the aggregate disk throughput.
- The Map tasks consume more CPU time than the Reduce tasks.

4. [7 points]: Why does the shuffle begin a long time before the Map phase has finished? Circle the best answer.

- There are more Map tasks ($M = 15,000$) than there are machines.
- There are more Reduce tasks ($R = 4000$) than there are machines.
- There are more Map tasks than there are Reduce tasks.
- The aggregate network throughput is smaller than the aggregate disk throughput.
- The Map tasks consume more CPU time than the Reduce tasks.

Name: _____

II Linearizability

These questions concern the material from Lecture 4, Consistency and Linearizability.

You have a service whose state is a single string, and that exposes two RPC operations to clients: one operation appends the RPC argument to the state, and the other RPC operation returns the current state. The timelines below indicate the start time, end time, argument string, and reply string for each client operation. Ax indicates an append operation with argument x , and Ry indicates a read operation to which the server replied y . The vertical bars indicate the start and end times of each operation (the times at which the client sends the request, and receives the reply). The service's state string starts out empty at the beginning of each history.

For example,

```
C1: |---Ax---|
C2:   |---Ay---|
C3:  |--Ryx--|
```

means that client $C1$ sent an append RPC with “ x ” as the argument, $C2$ sent an append RPC with “ y ” as the argument, and $C3$ read the state and received the reply “ yx ”.

Name: _____

Consider this history, in which the reply string sent to C4 has been omitted:

```
C1: |---Ax---|
C2:   |---Ay---|
C3:           |---Az---|
C4:           |--R?--|
```

5. [6 points]: Which values could C4's read yield that are consistent with linearizability?
Circle all of the correct answers.

- xzy
- yxz
- yzx
- xy
- xz
- yx
- zy

Name: _____

Now look at this history:

C1: |-----Ax-----|

C2: |---Ay---|

C3: |--Ry--|

C4: |----R?----|

6. [7 points]: Which values could C4's read yield that are consistent with linearizability?
Circle all of the correct answers.

- y
- x
- yx
- xy

Name: _____

III GFS and Raft

After reading the GFS paper (*The Google File System* by Ghemawat *et al.*) and the Raft paper (Ongaro and Ousterhout's *In Search of an Understandable Consensus Algorithm (Extended Version)*), Ben replaces the GFS master with a new coordinator that uses Raft. The Raft-based coordinator provides the same functions as before but replicates the log of operations using 3 Raft peers. All other parts of GFS stay the same.

Which of the following statements are true? (Circle all that apply)

7. [6 points]:

- A. The coordinator can continue operation in the presence of network partitions without any additional monitoring infrastructure, if one partition with peers is able to achieve a majority.
- B. The coordinator can continue operation correctly even if one of the 3 peers has failed (and there are no other failures).
- C. None of the above are true

Name: _____

Ben also considers using Raft for chunk replication. He runs many Raft clusters and has the GFS master assign chunks to a specific Raft cluster (i.e., each chunk is assigned to one Raft cluster, consisting of a leader and two followers). GFS clients submit write and append operations for a chunk to the leader of the Raft cluster for that chunk (i.e., Ben's design doesn't implement the separate data flow). The leader of the Raft cluster replicates write and append operation using the Raft library. All other parts of GFS (e.g., assigning leases to the leader, client caching locations of chunk servers, reading from the closest server, and so on) stay the same. (You can assume that chunk servers have enough disk space for operations to succeed.)

Which of the following statements are true? (Circle all that apply)

8. [7 points]:

- A. Unlike the old design, Ben's design can achieve linearizability for chunk operations.
- B. Unlike the old design, Ben's design can continue operation despite the failure of one chunk server.
- C. By using Raft, Ben's design allows clients to perform more mutating chunk operations per second than the old design.
- D. Raft's snapshots allow a chunk server to catch up in a few seconds if has been down for a long time (assuming the same network as in the GFS paper).
- E. None of the above are true

Name: _____

IV Raft

Consider the Raft paper (Ongaro and Ousterhout's *In Search of an Understandable Consensus Algorithm (Extended Version)*). Ben wonders what the impact of network behavior is on Raft's performance. Ben runs a Raft-replicated server that receives many client requests. If the network delivers AppendEntries RPCs in order, Ben's Raft implementation is fast (i.e., completes many client requests per second). But, if the network delivers AppendEntries frequently out of order, Ben's Raft implementation performs badly (i.e., completes fewer client requests per second). Using the rules in Figure 2 explain why this is the case.

9. [6 points]:

Name: _____

V Lab 3A-3C

Alyssa is implementing Raft as in Lab 3A-3C. She implements advancing the `commitIndex` at the leader (i.e., last bullet of Leaders in Fig 2) as follows:

```
func (rf *Raft) advanceCommit() {
    start := rf.commitIndex + 1
    if start < rf.log.start() { // on restart start could be 1
        start = rf.log.start()
    }
    for index := start; index <= rf.log.lastindex(); index++ {
        if rf.log.entry(index).Term != rf.currentTerm { // 5.4
            continue // ***
        }
        n := 1 // leader always matches
        for i := 0; i < len(rf.peers); i++ {
            if i != rf.me && rf.matchIndex[i] >= index {
                n += 1
            }
        }
        if n > len(rf.peers)/2 { // a majority?
            DPrintf("%v: Commit %v\n", rf.me, index)
            rf.commitIndex = index
        }
    }
}
```

Assume that all omitted parts of Alyssa's code are correct.

Name: _____

Ben argues that the line marked with “****” could be replaced by a break statement so that the loop terminates immediately.

10. [7 points]: Explain what could go wrong if one adopted Ben’s proposal; please include a specific sequence of events to illustrate your answer.

Name: _____

VI More lab 3A-3C

Alyssa is implementing Raft as in Lab 3A-3C. She implements the rule for conversion to follower in her `AppendEntries` RPC handler as shown below:

```
func (rf *Raft) convertToFollower(term int) {
    rf.state = Follower
    rf.votedFor = -1
    rf.currentTerm = term
    rf.persist()
}

func (rf *Raft) AppendEntries(args *AppendEntriesArgs,
                               reply *AppendEntriesReply) {
    rf.mu.Lock()
    defer rf.mu.Unlock()

    if args.Term >= rf.currentTerm {
        rf.convertToFollower(args.Term)
    }

    ...
}
```

Assume that all omitted parts of Alyssa's code are correct.

Name: _____

11. [6 points]: Describe a specific sequence of events that would cause Alyssa's implementation to break the safety guarantees provided by Raft.

Name: _____

VII ZooKeeper

Refer to *ZooKeeper: Wait-free coordination for Internet-scale systems* by Hunt, Konar, Junqueira, and Reed, and to the notes for Lecture 9.

The code fragments below are simplified versions of how something like GFS or MapReduce might use ZooKeeper to elect a coordinator, and for that coordinator to store state such as the assignments of GFS data to chunkservers.

Suppose server S1 executes the following code to become elected and to then store coordinator state in /A and /B. Initially, znode /coord-lock does not exist, znode /A starts out containing A0, and znode /B starts out containing B0.

```
s = openSession()
if create(s, "/coord-lock", data="S1", ephemeral=true) == true:
    setData(s, "/A", "A1", version=-1)
    setData(s, "/B", "B1", version=-1)
```

12. [7 points]: Briefly explain why, for coordinator election, it makes sense that /coord-lock should be an ephemeral znode rather than a regular znode.

Name: _____

S1's create() finishes and returns true to indicate success. But just after that, and before ZooKeeper has received S1's setData() requests, ZooKeeper decides that S1 has failed, and ZooKeeper terminates S1's session.

After ZooKeeper terminates S1's session, server S2 runs this to become coordinator:

```
s = openSession()
if create(s, "/coord-lock", data="S2", ephemeral=true) == true:
    setData(s, "/A", "A2", version=-1)
    setData(s, "/B", "B2", version=-1)
```

However, S1 is actually still alive, and it proceeds to send the two setData() requests, and they arrive at ZooKeeper.

Then client C1 reads /B and /A and sees B2 and A2, respectively.

Now a different client, C2, reads /B, and then reads /A. Both reads succeed.

13. [6 points]: Given the way ZooKeeper works, what can C2 observe? Circle all of the possible read results.

/B	/A

B0	A0
B0	A1
B0	A2
B2	A0
B2	A1

Name: _____

VIII Grove

In the `ApplyReadonly` function in Figure 7, Ben decides to delete the check for `s.waitForCommitted()`. The new code is as follows:

```
func (s *Server) ApplyReadonly(op) Result {
    s.mutex.Lock()

    if s.leaseExpiry > GetTimeRange().latest {
        e := s.epoch
        idx, res := s.stateLogger.LocalRead(op)
        s.mutex.Unlock()
        return res
    } else {
        s.mutex.Unlock()
        return ErrRetry
    }
}
```

14. [7 points]: Explain why this modification can result in non-linearizable reads.

Name: _____

IX Distributed Transactions

MouseGPT is designing a distributed transaction system using two-phase commit and two-phase locking, as discussed in Lecture 12 and Chapter 9 of the 6.033 reading. The goal is to provide serializable results. The question arises of what should happen if a participant computer crashes while in the PREPARED state for a transaction. MouseGPT thinks that all-or-nothing atomicity would be satisfied if such a transaction were completely forgotten. So MouseGPT designs the system so that if a participant computer crashes and restarts while it is in the PREPARED state for a transaction that it's part of, the recovery software on that computer un-does any local modifications the interrupted transaction might have performed and releases its locks, and sends a network message to each other participant and to the TC to tell them to undo any changes made by the transaction and to release its locks.

- 15. [6 points]:** Explain why MouseGPT's plan would cause the system to produce non-serializable (incorrect) results.

Name: _____

X 6.5840

16. [1 points]: Which lectures/papers should we definitely keep for future years?

- MapReduce
- RPC, Threads, and Go
- Linearizability
- GFS
- Raft
- ZooKeeper
- Q+A on Lab 3A/3B
- Grove
- Transactions

Name: _____

17. [1 points]: Which lectures/papers should we omit?

- MapReduce
- RPC, Threads, and Go
- Linearizability
- GFS
- Raft
- ZooKeeper
- Q+A on Lab 3A/3B
- Grove
- Transactions

Name: _____

18. [1 points]: What can we do to improve the course?

End of Exam I

Name: _____