



*Department of Electrical Engineering and Computer Science*

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

**6.824 Distributed System Engineering: Spring 2022**

**Exam II**

Write your name on this cover sheet. If you tear out any sheets, please write your name on them. You have 120 minutes to complete this exam.

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.

<b>I (7)</b>	<b>II (14)</b>	<b>III (14)</b>	<b>IV (14)</b>	<b>V (14)</b>	<b>VI (7)</b>	<b>VII (7)</b>	<b>VIII (14)</b>	<b>IX (7)</b>	<b>X (2)</b>	<b>Total (100)</b>

**Name:**

**Submission Website E-Mail Address:**

# I Frangipani

Consider the paper *Frangipani: A Scalable Distributed File System*, by Thekkath *et al.* Four users,  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ , each with their own workstation, share a Frangipani file system. The file system appears under the pathname  $/f$ . The file system starts out containing no files and just the empty directory  $/f$ .

At about the same time,  $\alpha$  runs

```
cp /dev/null /f/a
```

and  $\beta$  runs

```
cp /f/a /f/b
```

and  $\gamma$  runs

```
cp /f/b /f/g
```

`cp` is a command that copies one file to another. If the source file doesn't exist, `cp` prints an error message and does *not* create the destination file. `/dev/null` is an empty file, so the first command is a way to create a new file (`/f/a`) containing nothing.

All the above commands complete and indicate that they were successful. There is no other activity, and no failures up to this point.

Now the workstations of  $\alpha$ ,  $\beta$ , and  $\gamma$  crash!

**1. [7 points]:** After the crashes, user  $\delta$  looks at the Frangipani file system,  $/f$ , from her workstation. What might  $\delta$  see as the *complete* output of a directory listing of  $/f$ ? Circle all correct answers.

- A. No files.
- B. `/f/a`
- C. `/f/b`
- D. `/f/g`
- E. `/f/a /f/b`

## II Spanner

2. [7 points]: Sections 4.1.4 and 4.2.2 of *Spanner: Google's Globally-Distributed Database* say that in many cases Spanner assigns a read-only transaction a timestamp of `TT.now().latest`. What useful property does this choice of timestamp have? Circle the single best answer.

- A. This timestamp is guaranteed to be less than  $t_{safe}$  at every replica.
- B. This timestamp is guaranteed to be less than the timestamp of any concurrent read-write transaction.
- C. This timestamp is guaranteed to be greater than the timestamp of any completed read-/write transaction.
- D. None of the above.

3. [7 points]: For a read-only transaction that reads multiple keys stored on different Paxos groups, why wouldn't it be correct to read the most up-to-date version of each of the keys? Please briefly describe an execution in which the results would be incorrect.

### III FaRM

Consider the FaRM execute/commit protocol in Figure 4 of *No compromises: distributed transactions with consistency, availability, and performance*.

The following two transactions run on different computers that share a FaRM storage system:

T1:

```
txCreate()           // start the transaction
ox = txRead(x_oid)  // read object x (an integer)
txWrite(x_oid, ox + 1) // write object x
oy = txRead(y_oid)  // read object y (an integer)
oy = compute using ox and oy for a while
txWrite(y_oid, oy)  // write object y
txCommit()          // start Figure 4's commit phase
```

T2:

```
txCreate()
ox = txRead(x_oid)
oy = txRead(y_oid)
ok = txCommit()
if ok:
    print ox, oy
```

`txRead()` uses RDMA to read an object with a given object identifier into local memory. `ox` and `oy` are local variables. `txWrite()` buffers the written data, which is written to FaRM during the commit phase if the transaction commits. `txCommit()` returns true if the transaction committed successfully, and false if it aborted.

Suppose T1 is in its Execute phase, and has just read object x. At this point, T2 starts executing, on a different computer. Nothing else is happening. There are no failures.

**4. [7 points]:** Could T1 cause T2 to abort? Explain briefly how this could happen, or why it can't happen.

**5. [7 points]:** Is it possible for T2 to successfully commit? Explain briefly how this could happen, or why it can't happen.

## IV Spark

The paper *Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing* by Zaharia *et al.* mentions a distinction between wide and narrow dependencies.

**6. [7 points]:** Some kinds of dependencies, such as join, are sometimes wide and sometimes narrow. Join is a transformation that brings together records from its two input RDDs that share the same key (so a join's inputs must be RDDs of key/value pairs). Circle the single answer that best explains when join is wide versus narrow.

- A. Join is narrow when executed lazily, and wide when forced to materialize by an action such as `collect()`.
- B. Join is narrow when its two input RDDs resulted from narrow dependencies, wide otherwise.
- C. Join is narrow when both preceding RDDs are partitioned by key in the same way and the join is joining on that key, wide otherwise.
- D. Join is narrow when the preceding RDD is persistent, wide otherwise.
- E. None of the above.

**7. [7 points]:** Circle the true statements (there may be more than one).

- A. A wide dependency requires data to be moved across the network, while a narrow dependency does not.
- B. A narrow dependency can be executed lazily, while a wide dependency must be executed immediately.
- C. A narrow dependency requires input RDD data to be key/value pairs, while a wide dependency's input can be arbitrary records.
- D. A narrow dependency is likely to execute more slowly than a wide dependency.
- E. If a worker fails, narrow RDD partitions for which it was responsible may have to be re-computed starting with the original inputs.

## V Memcached at Facebook

Figure 1 of *Scaling Memcache at Facebook* by Nishtala *et al.* illustrates a straightforward way to use memcache as a cache for a database. The left-hand part of Figure 1 shows how an application reads data, and the right-hand part shows how an application writes data. For these questions, assume the application caches database records in memcache, with the same keys that the database uses.

**8. [7 points]:** Consider the scheme in Figure 1 **without leases** and without any of the paper's other improvements (without McSqueal, without McRouter). Why does Figure 1's write path do the memcache delete *after* the database UPDATE, rather than before? Circle the single best answer.

- A. This order guarantees that a read concurrent with a write of the same key won't cause stale data to be cached in memcache.
- B. This order prevents a read from seeing the old value after the UPDATE completes.
- C. Neither of the above.

**9. [7 points]:** Now consider a design consisting of Figure 1 **with leases** (but without any of the paper's other mechanisms). In which situation(s) does the lease mechanism decrease the likelihood of stale data being cached in memcached indefinitely, compared to Figure 1 alone? In all of these scenarios, the read's get misses in memcache. Circle all of the correct answers.

- A. When there are two concurrent reads of the same record (but no writes).
- B. When the two steps of a write occur between steps 1 and 2 of a read of the same record.
- C. When the two steps of a write occur between steps 2 and 3 of a read of the same record.
- D. When the three steps of a read occur between steps 1 and 2 of a write to the same record.
- E. None of the above.

## VI COPS

Alyssa stores data in COPS (*Don't Settle for Eventual: Scalable Causal Consistency for Wide-Area Storage with COPS* by Lloyd *et al.*). The COPS deployment has three replicas, in three data-centers. Alyssa runs code on three clients (C1, C2, and C3), each in a different data-center along with one of the COPS replicas. Alyssa is using COPS (not COPS-GT).

Alyssa's data consists of two key/value pairs: I1, whose value is a string, and I2, whose value is an integer. I1's value starts out as the empty string, and I2's value starts out as zero. Both key/value pairs start out identical at the three COPS replicas.

At about the same time, each of Alyssa's three clients runs some code that calls the COPS client library put and get functions. There's no other activity, and no failures.

C1 runs this:

```
put("I1", get("I1") + "X") // append X to I1
put("I2", get("I2") + 1)   // increment I2
```

C2 runs this:

```
put("I1", get("I1") + "X") // append X to I1
put("I2", get("I2") + 1)   // increment I2
```

C3 runs this:

```
print "I2=", get("I2") // read I2 first
print "I1=", get("I1") // then I1
```

**10. [7 points]:** What output could C3 print? Circle all of the correct answers.

- A. I2=2 I1="XX"
- B. I2=2 I1="X"
- C. I2=2 I1=""
- D. I2=1 I1="XX"
- E. I2=1 I1=""
- F. I2=0 I1="XX"



## VII SUNDR

Three client computers (C0, C1, and C2) share files on a SUNDR server. (*Secure Untrusted Data Repository (SUNDR)*, by Li *et al.*). The SUNDR server starts out containing three files, x, y, and z, and their contents all start out as 0.

Client C0 asks the SUNDR server to perform the following three writes, one at a time, waiting for each to complete before proceeding to the next:

```
C0: write "1" to file x
C0: write "2" to file y
C0: write "3" to file z
```

The SUNDR server returns a success indication for each of these writes, and C0 does not detect any problem. After C0 has completed its writes, C1 asks the SUNDR server to perform the following three writes, one at a time, waiting for each to complete before proceeding to the next:

```
C1: write "1" to file x
C1: write "2" to file y
C1: write "3" to file z
```

Again, the SUNDR server returns a success indication for each C1's writes, and C1 does not detect any problem. After C1 has completed its writes, C2 asks the SUNDR server to perform the following three reads, one at a time, waiting for each to complete before proceeding to the next:

```
C2: read file x
C2: read file y
C2: read file z
```

The SUNDR server returns a success indication for each of these reads, and C2 does not detect any problem.

The three clients only communicate with the SUNDR server; they don't talk to each other. You can answer this question with reference to the paper's straw-man protocol (Section 3.1), though that should not affect the answers.

**11. [7 points]:** Which of the following could C2 see from its three reads?

**A.**  $x=1$   $y=2$   $z=3$

**B.**  $x=0$   $y=2$   $z=3$

**C.**  $x=0$   $y=0$   $z=0$

**D.**  $x=1$   $y=0$   $z=0$

**E.**  $x=1$   $y=0$   $z=3$

**F.**  $x=1$   $y=1$   $z=1$

## VIII Bitcoin and Ethereum

Ben runs an online tulip shop where he accepts Bitcoin as payment, 1 BTC for 1 tulip (these are premium tulips). Upon seeing a payment transaction appear in the blockchain, Ben waits for 6 additional blocks before treating the payment as valid and shipping tulips to the buyer by overnight Fedex.

- 12. [7 points]:** Suppose an attacker owns 1 BTC, and the attacker gains control of 80% of the Bitcoin network hash power for a period of 2 hours. Describe how this attacker could cheat Ben's store, such that on the next day, the attacker still has their 1 BTC and also has a tulip.

Alyssa runs a virtual tulip shop implemented as an Ethereum smart contract. Her shop sells tulip "tokens" for 1 ETH (ETH is the Ethereum currency). A tulip token is a purely virtual item, whose ownership is represented in the state of Alyssa's smart contract; no physical goods are involved. Here's what Alyssa's Ethereum tulip smart contract state and external method look like:

```
// track how many tulip tokens each ETH address holds
private mapping (address => uint256) tulips;

function buy() external payable {
    require(msg.value >= 1 eth); // sender has to transfer >= 1 ETH
    tulips[msg.sender]++; // sender gains a tulip token
}
```

- 13. [7 points]:** Suppose an attacker owns 1 ETH, and the attacker gains control of 80% of the Ethereum network hash power for a period of 2 hours. Could the attacker cheat Alyssa's store, such that on the next day, the attacker still has their 1 ETH and also has a tulip token? Why or why not?

## IX Lab 3

Alyssa is working on reducing the read latency of her Lab 3 by allowing either a leader or a follower to respond to a Get() RPC, and to do so without putting the operation in the Raft log. Her plan is as follows:

- A. Modify the Put() and PutAppend() RPC responses to return the Raft log index where the operation committed.
- B. Modify the client Clerk module to track the latest log index received in a Put() or PutAppend() RPC response, and pass the index as an additional lastPutIndex argument to the Get() RPC request.
- C. Modify the server code to execute Get() requests with a given lastPutIndex as follows, on either the leader or a follower:
  - (a) If the local replica has finished processing at least up through index lastPutIndex, then read the data from the local replica's state and return immediately.
  - (b) Otherwise, wait until the command at lastPutIndex has been processed, and then read the data from the local replica's state.

**14. [7 points]:** Alyssa claims that her design will order her Get() requests after previous Put()/PutAppend() requests, and this will prevent stale reads. Is Alyssa's modified design correct — does her design still implement a linearizable key-value store? If yes, explain why it's correct. If not, give a scenario in which Alyssa's design results in an execution that violates linearizability.

## **X 6.824**

**15. [1 points]:** Which lectures/papers should we **omit** in future years?

- Distributed Transactions
- Frangipani
- Spanner
- FaRM
- Spark
- Memcached at Facebook
- COPS
- SUNDR
- Bitcoin
- Blockstack
- Ethereum

**16. [1 points]:** Do you have any feedback for us about 6.824?

**End of Exam II**