

Notes:

There is some overlap in what I wrote below and what you had written. I wrote it partially to try to get it straight in my mind. We should probably reorganize / merge some things (the discussion of replication for instance). I'm also not sure we need to do replication for project 3, only for project 4, so we could put of implementing this and note our plan here.

1 Commands:

Our system has two kinds of nodes: managers and clients.

Managers currently don't support any commands.

Clients support the following commands:

Configuration Commands: `manager`, `manageris <addr>`, `handshake <addr>`

FS Commands: `create <filename>`, `delete <filename>`, `get <filename>`,
`put <filename> <contents>`, `append <filename> <contents>`

TX Commands: `txstart`, `txcommit`, `txabort`

Dev Debug Commands: `debug`, `noop <addr>`

All operations performed outside of transactions are semantically implicitly wrapped inside of a `txstart` and `txcommit` by the manager. However, in our actual implementation, different code runs if the client is or isn't currently performing a tx for performance reasons.

2 File System Semantics:

A file is defined to exist and is accessible iff some node in the system created the file outside of a tx, or inside of a tx that has been committed. This file will cease to be available if it is deleted outside of a tx or deleted by a committed tx.

If a file exists, then it is guaranteed that the manager has some version of the file in its persistent storage (although, it is not guaranteed to be the newest version). As with our previous project, we allow one client at most to have ReadWrite access to a file.

Since we are not using write-through cache coherence, we expect clients to maintain responsibility for files they have RW access on. We expect that clients will use replicas to appropriately clone data so that the data can be recovered if the client goes down. This also means that the manager will not request the updated version of these files until another client requests them. At this point, the manager will ask the client to send its changed version of the files forward and to revoke its local permissions on the file(s) (and the manager will revoke the client's permissions on the file(s) as well).

In the event of a client failing while that client still has RW access on a file(s), the manager will transfer ownership to the client's replica and then request the data again.

If no client has ReadWrite access to a file, any number of clients can have ReadOnly access to that file. The manager explicitly invalidates these permissions before granting anyone else ReadWrite access to that file.

3 Serialization:

We employ client-side file-level locking to prevent a client from requesting permission to the same file twice in a row. For example, say a client doesn't have ReadWrite access to the file test.txt and the client receives the commands:

```
put test.txt I'm about to delete this file!
delete test.txt
```

Since the client gains RW access on test.txt after the first command, there is no need to contact the manager to perform the second command, so we queue the second operation until the first completes.

The manager, in turn, assumes the client will queue commands until they are ready for the next request to be serviced. If a client requests access to the same file twice, the manager will grant them the appropriate level of access.

4 Transaction Scheme:

We decided to use two-phase locking manager-side to ensure serializable transactions for the following reasons:

Framework setup Our framework is already suited towards two-phase locking. We already have a locking scheme in-place for cache coherency (project 2), and so the process of locking files was simply expanded to lock files for the duration of a transaction, instead of for the duration of a request.

Difficulties with optimistic concurrency We thought about using optimistic concurrency, however this requires the server to validate all requests at the end of a transaction. This would require a log on the server recording transactions. Further, this log would have to timestamp all transactions in order to decide whether two transactions conflicted or not. Lastly, this would also required the server to be sent a copy of the client's transaction log with every commit, which we thought was unnecessary.

Versioning Further, using optimistic concurrency would require some form of file versioning, in order to support rollback. We thought this approach required more space allocation than was necessary or ideal.

5 Deadlock:

One potential problem with 2PL is the possibility of deadlock.

We require clients to operate on files within a transaction in filename order to avoid deadlock. So, if a client wants to perform the following operations:

```
Get g => x
Put x => f
```

They have to actually perform these operations:

```
Append "" f
Get g => x
Put x => f
```

Following this procedure, clients are guaranteed not to dead lock while performing a transaction.

The server currently makes no guarantees to clients if they do not perform requests in filename order. Clients can always abort themselves, which results in all locks freeing eventually. It would not be difficult for the manager to ensure clients lock files in order, but we haven't implemented it yet.

6 Transaction Semantics:

A client is considered to be transacting until it receives a TX_SUCCESS or TX_FAILURE from the manager. The client will queue all commands it is given until it receives a tx response from the manager for its outstanding transaction.

The manager will not wait for an ack on the TX_SUCCESS or TX_FAILURE packet before revoking that client's permissions. The RIO message layer ensures that a delayed or dropped packet won't cause a problem: if a TX_SUCCESS packet is delayed and the manager requests a file from the client, the TX_SUCCESS will be acted upon before any request from the manager is serviced.

7 TFS Log Entry Format:

```
<client_address>
<Operation type>
<filename>
<contents_line_count>
<contents>
```

contents_line_count is -1 for operations that don't have contents there is a line separator after each *< entry >* including the contents tx ops don't have a filename or anything after

8 Failure Handling:

Server Failures Currently, clients block on server failures. This is within the specification of the assignment, and so we did not implement any handling of this scenario. This will be relaxed by PAXOS in assignment 4.

Client Failures While a client is in the middle of a transaction, they are periodically sent pings by the server (heartbeat pings). If the client ceases responding to this heartbeat, then after a set number of rounds the server will assume the client went down (this functionality was largely borrowed from project 1, where servers were required to deal with client failures). If the server detects that a client went down, it will immediately release all locks this client had on any files and abort their transaction.

The server will then immediately change ownership of any files that client had ownership of to its replica. If there were pending permission requests for this file, then the server will immediately forward that request to the appropriate replica. TODO: HIGH: NOTE: I think we might want to implement this w/ RPC or a different message type for implementation simplicity.

Transaction Failures Since all transaction commands are written to a log before actually being committed, a transaction failure results in no changes to the local or remote file system occurring. Instead, the client may attempt to redo the transaction based on the type of failure that occurred, or it may simply decide to abort that transaction altogether, at which point the user would have to redo the transaction.

Command Failures Command failures result in an automatic transaction failure for our framework. We chose to implement this feature as opposed to the alternative (sending errors but proceeding with the transaction) because we believe that users could end up in an undesirable state if they commit a transaction that succeeds on some commands but not others. Rather than allow the client to end up in an undesirable state, we decided to abort the transaction and force the user to recommit a new transaction without the offending command.

9 Replication:

For project 4, we intend to implement the following replication scheme:

In the middle of transactions, in order to avoid potential situations where the most current copy of the file is lost, we decided to implement replication. Whenever a client changes a file locally, it will also replicate this action via RPC on another client (we chose a simple replication scheme, where client 1 replicates on client 2, ... client k replicates on client k+1, ... client N replicates on client 1). This ensures that the latest copy of a file will never be lost, even if a client fails (see below for further discussion on that topic).

10 Running Test Scripts:

Test scripts should be called as follows (a sample test script is given as an example):

```
./clean.sh; ./compile.sh; ./execute.pl -f 0 -n Client -s -c ./simulator_scripts/test_tx_3
```