<div align="center">

CSE 490h – Project 2: Cache Coherence Writeup
**Wayne Gerard - wayger**
**Zachary Stein - steinz**
January 27, 2011

</div>

**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.

**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 implicity wrapped inside of a txstart and txcommit by the manager. However, in our actual implementaion, different code runs if the client is or isn't currently performing a tx for performance reasons.

**File System Semantics:**

A file exists and will be accessible to any node in the system iff it has been created by some node in the system outside of a tx or during a committed tx and it hasn't been deleted outside of a tx or during a committed tx since it was created.

Any file existing in the system is guaranteed to have some version in the manager's persistent storage (although this might not be the newest version). For a given file, at any point in time at most one client can have ReadWrite access to that file. If some client has RW access to a file, it is that client's responsibility to make sure that its replicas log all changes the owning client makes to that file before the client reports changes completed either to the manager or to the command-issuer. We employ write-back cache coherence, so the manager will not get the updated version of the file until someone else requests it. When this happens, the manager tells the client to report its changes and lose RW. If the client fails to respond to the manager's request, the manager asks one of the client's replicas for the file instead.

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.

If a client is slow or network failure occurs, the manager might be forced to give someone ReadWrite access to a file before receiving all invalidation confirmations from the clients who previously had ReadOnly access to that file. This means that after reconnection, two clients might

<div align="center">

1

</div>

think they have ReadWrite, so the manager always verifies that clients trying to perform mutations on files have RW access to those files.

**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 fist completes.

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.

imistic 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.

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

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.

**Replication:**

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).

**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). 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 simplicity.

**TFS Log Entry Format:**

¡client$_a$ddress >< Operationtype >< filename >< contents$_l$ine$_c$ount >< contents >

contents$_l$ine$_c$ountis$-1$foroperationsthatdon'thavecontentsthereisalineseparatoraftereach < entry > includingthecontentstxopsdon'thaveafilenameoranythingafter