# Team Members:

Neel Yadav (34064006)
Aneri Rana (33157215)

# Design approach for Zookeeper:

In the init method, we create a new zookeeper session and cassandra cluster. To use zookeeper's coordination services for our fault tolerant replicated database we use a combination of sequential nodes and watcher in the following manner:

1. In handleMessageFromClient(), any server that receives a new request from the client will add a new znode of persistent sequential type with the request value under '/requests' znode.
2. Zookeeper coordinates the auto-increment of request numbers for us as multiple servers try to add their requests.
3. Every server registers a watcher on children of '/requests' during initialization, so that each server is notified when any server receives and adds a request.
4. The watcher triggers rollForward() method, which takes the last executed request number for the server and then executes all the remaining requests in the /requests starting from the last executed one (this is done so that if any request triggers were missed they will still be executed in a sequential manner).

**Checkpoint and recovery**
We decided to select a checkpoint threshold of 50 requests to find a balance between decent frequency of backups while not creating too much backup overhead. When the last executed request number reaches a multiple of 50, the checkpoint() method is called. The checkpoint() method triggers the creation of a snapshot by calling takeSnapshot. The takeSnapshot method creates a snapshot of the Cassandra keyspace by fetching all rows from the table and serializing the data to a .ser file.

We do crash recovery during initialization. If a checkpoint exists for a server, then we call the restoreSnapshot() and rollForward() methods to bring the server back upto speed if it is recovering from a crash. The restoreSnapshot method reads the latest snapshot file and restores the data back into the keyspace table by creating insert queries for each id. If there is no checkpoint stored but there are pending requests logged, then we execute all the requests from /requests to get the server upto speed.
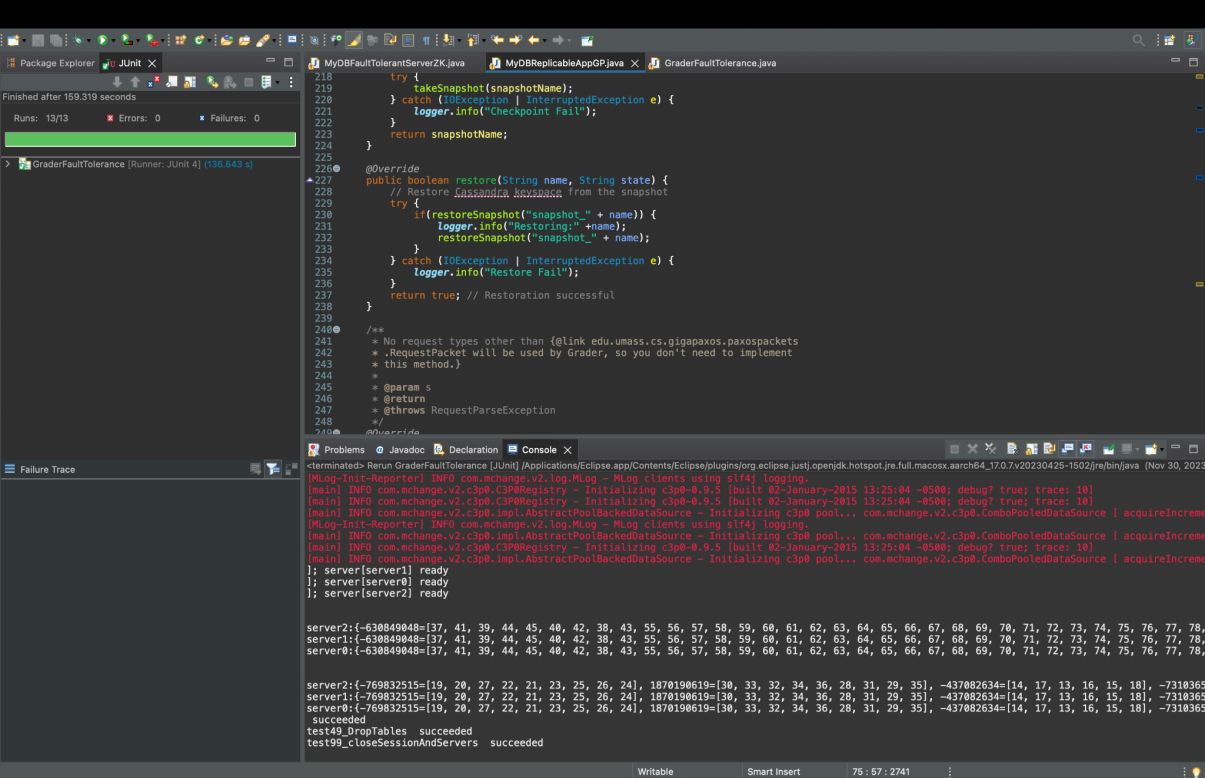
Additionally, we also store the last executed request number in the snapshot filename with format "checkpoint_[reqNum].ser". This is used during recovery so that any pending requests are executed after the last executed request number only.

**Trim Logs**
Everytime we have reached the checkpoint threshold of 50 requests, we call trimLogs() after taking the snapshot. If the log capacity has reached 400 then the oldest 50 requests would be pruned to make sure the logs don't grow infinitely long.

# Design approach for Gigapaxos:

We extend the replicable API of gigapaxos and implement the three functions, execute, checkpoint and restore to successfully implement a consistent cassandra database. We were able to pass all the test cases.



The execute method handles the execution of incoming requests. It extracts the CQL query from the request and executes it using the Cassandra Session. The takeSnapshot method creates a snapshot of the Cassandra keyspace by taking all rows from the table in the keyspace and serializing the data to a ser file. The restoreSnapshot method reads a snapshot file and restores the data back into the table in the table keyspace by creating insert queries for each id. The checkpoint method triggers the creation of a snapshot by calling takeSnapshot and returns the snapshot name. The restore method restores the keyspace from a given snapshot name using restoreSnapshot.