

# CS60002 Distributed Systems

## *Assignment 2: Consensus Module using Raft*

Maitrey Govind Ranade (18CS30026)

Kothapalli Dileep

Arpit Das

Gawai Laukik

### ***Part 1: The ATM machine***

We use PySyncObj available at <https://github.com/bakwc/PySyncObj> as the choice of our RAFT library. This library allows us to create RAFT objects in a python program. These objects then can be synced to other running python programs. For creating an ATM system, we create a single RAFT object that maintains a balance corresponding to accounts. The object is updated when the following defined functions are called on the raft object, which update the in-memory dictionary of python maintaining a balance corresponding to each account by reading the RAFT logs:

- **createAccount(account):** append to RAFT Logs (indicated by @replicated decorator).  
Adds a key in data dictionary equal to account with zero balance
- **deposit(account, value):** append to RAFT Logs (indicated by @replicated decorator).  
Deposit the amount = value to account
- **withdraw(account, value):** append to RAFT Logs (indicated by @replicated decorator), Withdraw amount = value from account
- **transfer(account1, account2, value):** append to RAFT Logs (indicated by @replicated decorator): Deposit amount = value from account1 to account2
- **inquire(account):** Read-only function, get the balance of account.

Part 1 can be tested as follows:

```
Python tests/atm_test.py
```

It will open 3 ATMs on separate terminals. you can enter commands in the terminal to withdraw, deposit or transfer in the opened terminal

## Part 2: Fault Tolerant Broker

### Overview

The majority of the assignment is same as assignment 2. We just update the broker server to instantiate RAFT objects for each partition they handle that sync with the corresponding replica partitions on other brokers using our RAFT library of choice. We also change the broker manager appropriately to maintain the broker replica information. We now give an updated description of each module of the source code.

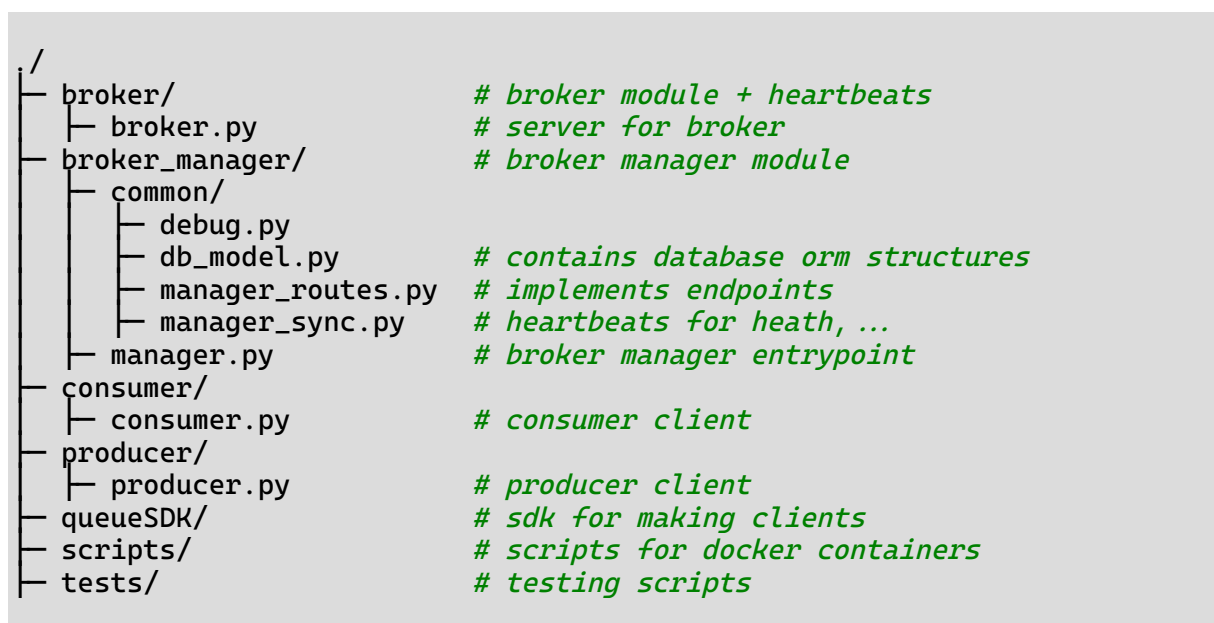
We design all the specified modules in the assignment i.e. :

- Broker
- Broker Manager (read/write)

as Flask applications. Flask framework provides an easy way to build http endpoints, enabling **asynchronous** and **simultaneous** handling of the designed endpoints for various tasks by using **threading**.

For persistent data storage required by any process, it uses postgresSQL, with **write-ahead logs (WAL)** enabled to ensure data protection against crashes.

### Source Structure



### Broker

The broker module is a flask application contained in the folder `broker/broker.py`. It supports the following endpoints

- **/store\_message:** Endpoint for receiving a message from the primary (write) broker manager. It receives message and writes it into the database.
- **/retrieve\_messages:** Retrieve multiple messages with IDs greater than a given offset.
- **/consume:** Endpoint for consumer request as directed by the read manager

Design Choices:

- The broker maintains an in-memory raft objects for each partition. The RAFT library we used is **modified by us** to extract the RAFT logs and commit the logs to the postgresql database running on. The manger for the RAFT objects runs on a single port 4001 as mentioned in the assignment statement. The PySyncObject has issues logs in the form of of COMMANDS. We store these commands on our postgres database so tat these can be recovered in the event of a failure.
- The redundancy factor is kept to 3 replicas per partition, which means all the partitions are guaranteed to remain intact and functional in the event of a single failure.
- The since the brokers have the actual partition messages, they have RAFT object to maintain the consumer offset for each partition with is synced across replicas. This has some performance advantages as consumer requests can be directly redirected to a broker by the manager server.

The brokers periodically send **heartbeats** to write broker manager. They include their IPs and ports in each beat. So the heartbeats of the brokers also double as a **discovery service** for the broker, serving two purposes at once. (health checks and discovery).

## Broker Manager

The broker manager module is a flask application contained in in the file `broker_manager/manager.py`. It can be run in a **read mode** (reads messages from the brokers for consumers) or a **write mode** (writes messages by simply passing an argument.

Each broker manager (read/write) maintains the following tables for **metadata** storage:

- **Producer** maps id, topic, (optionally partition) and health
- **Consumer** maps id, topic, (optionally partition) and health
- **Topic:** maps topic id and name
- **Partition:** maps id, topic id, broker id, replica
- **Broker:** contains broker ids and their ips and ports. Also maintains health status

- **Replica:** contains ips and ports and health status of replicas.

Since we use PostgreSQL as our choice of storage at the brokers. Making more partitions will not improve message retrieval performance (due to a fixed schema). Our design supports any number of partitions as per assignment. But as a design choice, every topic has a default (implicit) partition on every broker (with id -1). If a producer does not specify a partition ID, its messages go to any random broker in the implicit partition. On the other hand if a consumer does not specify a partition ID, it consumes messages from all possible partition IDs.

We now describe the endpoints supported by broker manager in our design, with the relevant design choices:

Endpoints Open in Both modes:

- **[GET] /topics:** Get a list of created topics.
- **/topics/partitions:** Get a list of partitions of a given topic name (returns broker info of partitions)
- **/brokers:** returns healthy brokers.
- **/replicas:** List all the known healthy replicas of broker manager in the system.

Endpoints Open in Write mode:

- **/brokers/create/partition:** Creates partition for a given topic name and broker ID.
- **/brokers/remove:** Removes a broker with a given ID
- **/brokers/heartbeat:** Receives heartbeat from a broker (with broker IP and port), acts as a broker discovery service for write broker manager.
- **[POST] /topics:** Creates a new topic with a given topic name and returns its ID
- **/partitions:** Creates a partition given a topic name and a broker ID on the broker.
- **/producer/register:** Registers a producer for a given topic name and partition ID. If Partition ID is not given or -1, it registers for the common implicit partition for every topic on every broker.
- **/producer/produce:** Takes a message by producer and requests a broker to store it. If a partition ID is supplied, it will communicate and store to **any replica** at random the partition. Otherwise it will choose a partition **uniform at random** from healthy brokers and put the message in its implicit partition for uniform distribution, leading to uniform load distribution by dequeue.

- **/consumer/health\_poll:** (Only opened by primary) Endpoint for inter-broker manager communication. The consumers poll the read manager, the read manager redirects the polling to this endpoint for health metadata update for consumers.
- **/metadata/sync:** Read only managers poll this endpoint on primary to sync metadata with it.

Endpoints Open in Read mode:

- **/consumer/register:** Registers a consumer for a given topic name and partition ID. If Partition ID is not given or -1, it registers for all the partitions of the given topic name.
- **/consumer/consume:** Endpoint for consumer to dequeue. If a partition ID is specified, it will consume from the any replica of the specified partition. Otherwise, it will choose a partition **uniformly at random** and dequeue the latest message for the given topic in that particular broker. This is done to maintain uniform load balancing. Updates health info using /consumer/health\_poll.

In addition to all these endpoints, the primary broker manager runs a thread to periodically update health status if brokers, producers, consumers and replicas.

Replicas run a thread periodically to sync their metadata with the primary.

## Other Components

Other components such as SDK for queue (in folder queueSDK/) and producer and consumer processes are updated from the previous assignment according to new newly designed endpoints.

## Testing

- We do our testing using docker containers.
- Each node (broker, write broker manager, read broker manager) is launched inside its own docker container (with their own database).
- All the nodes communicate through the default docker subnet (172.17.0.0/16).
- IP of write manger is fixed to 172.17.0.2, as write manger also acts as a service discovery service for all other nodes spawned later on.
- Producer and consumer clients are launched directly on the host system.
- Node crashes are simulated by using the `docker kill` command.

The `scripts/` directory contains helper scripts for testing and launching containers.

- `create_image.sh`: creates the base docker image for all nodes
- `launch_container.sh <CONTAINER-NAME> <PROCESS-FILE>.py`: Launches a fresh container for a node with flask application `<PROCESS-FILE>.py`.
- `boot_container.sh <CONTAINER-NAME> <PROCESS-FILE>.py`: Starts a killed/stopped container for a node with flask application `<PROCESS-FILE>.py`.

To quickly test everything, just run the created testcase:

```
./tests/test.sh
```

## Discussion

Although the brokers are now immune to crash faults and messages won't be lost if a broker crashes, the same can not be said for write manger. The information of the write manager is replicated across read mangers. But if the write manager fails, the producers can not produce and new producers and consumers can not be registered.

The complete solution for the bonus part involves having a raft object that maintains the metatdata consistently cross all read/write managers. Now even if the write manager fails, the producer can continue producing to any manager, since all managers are consistent by using RAFT. We just did not attempt the implementation for the bonus part due to lack of time.