

CS60002: Distributed Systems

Assignment 2: Implementing a Scalable Database with Sharding

Date: Feb 06th, 2024

Target Deadline (For sharing the git repo): March 15, 2024

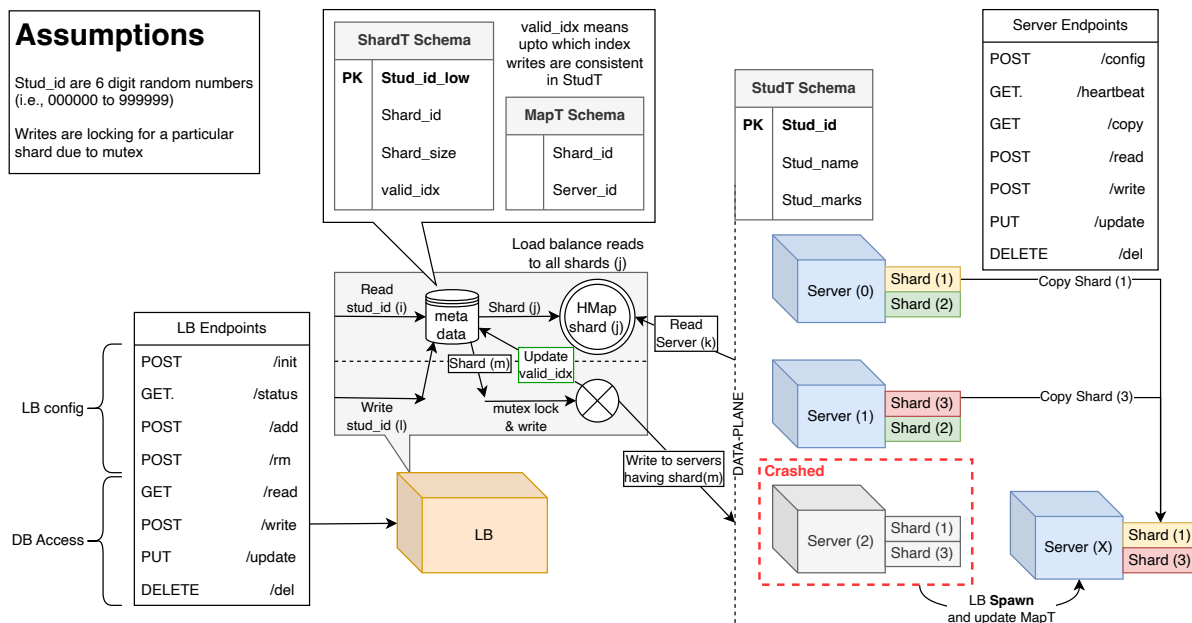


Fig. 1: System Diagram of a Sharded Database

I. OVERVIEW

In this assignment, you have to implement a sharded database that stores only one table *StudT* in multiple shards distributed across several server containers. This is an incremental project so that you can reuse the codebase from the first assignment. A system diagram of the sharded database is shown in Fig. 1. Here, shards are subparts of the database that only manage a limited number of entries (i.e., *shard_size* as shown in the diagram). Shards can be replicated across multiple server containers to enable parallel read capabilities. For this assignment, we assume that write requests are blocking for a particular shard. Thus, if two write requests are scheduled simultaneously on shard (i), one of them will wait for the other to complete. However, Parallel writing to different shards, for instance, shard (i) and shard (j), is possible. The system's current design provides scaling in two ways: (i) Read speed with more shard replicas and (ii) Database size with more shards and servers.

A. Coding Environment

- **OS:** Ubuntu 20.04 LTS or above
- **Docker:** Version 20.10.23 or above
- **Languages:** C++, Python (preferable), Java, or any other language of your choice

B. Submission Details

- Write clean and well-documented code.
- Add README file and mention the design choices, assumptions, testing, and performance analysis.
- Add Makefile to deploy and run your code.

Your Code should be version-controlled using Git, and the GitHub repository link must be shared before the deadline. Please note that the contribution of each group member is essential to learn from these assignments. Thus, we will inspect the commit logs to award marks to individual group members. An example submission from previous course session can be found at <https://github.com/prasenj52282/shardQ> for your reference.

II. TASK1: SERVER

The server containers handles shards of the *StudT* (**Stud_id: Number, Stud_name: String, Stud_marks: Number**) table. Each server container can manage a different set of shard replicas of the distributed database, as shown in Fig. 1. The endpoints to handle requests to a specific server are described as follows:

- 1) **Endpoint (/config, method=POST):** This endpoint initializes the shard tables in the server database after the container is loaded. The shards are configured according to the request payload. An example request-response pair is shown below.

```
1 Payload Json= {
2   "schema": { "columns": ["Stud_id", "Stud_name", "Stud_marks"],
3               "dtypes": ["Number", "String", "String"] }
4   "shards": ["sh1", "sh2"]
5 }
```

```

6 Response Json = {
7     "message" : "Server0:sh1, Server0:sh2 configured",
8     "status" : "success"
9 },
10 Response Code = 200

```

- 2) **Endpoint (/heartbeat, method=GET):** This endpoint sends heartbeat responses upon request. The heartbeat endpoint is further used to identify failures in the set of server containers maintained in the distributed database. Therefore, you could send an empty response with a valid response code.

```

1 Response [EMPTY],
2 Response Code = 200

```

- 3) **Endpoint (/copy, method=GET):** This endpoint returns all data entries corresponding to one shard table in the server container. Copy endpoint is further used to populate shard tables from replicas in case a particular server container fails, as shown in Fig. 1. An example request-response pair is shown below.

```

1 Payload Json= {
2     "shards":["sh1","sh2"]
3 }
4 Response Json = {
5     "sh1" : [{ "Stud_id":1232, "Stud_name":ABC, "Stud_marks":25},
6               { "Stud_id":1234, "Stud_name":DEF, "Stud_marks":28},
7               ....],
8     "sh2" : [{ "Stud_id":2255, "Stud_name":GHI, "Stud_marks":27},
9               { "Stud_id":2535, "Stud_name":JKL, "Stud_marks":23},
10              ....],
11     "status" : "success"
12 },
13 Response Code = 200

```

- 4) **Endpoint (/read, method=POST):** This endpoint reads data entries from a shard in a particular server container. The endpoint expects a range of *Stud_ids* (low, high) and *Shard_id* that has to be read from the server container.

```

1 Payload Json= {
2     "shard":"sh2",
3     "Stud_id": { "low":2235, "high":2555}
4 }
5 Response Json = {
6     "data" : [{ "Stud_id":2535, "Stud_name":JKL, "Stud_marks":23},
7               { "Stud_id":2536, "Stud_name":MNO, "Stud_marks":22},
8               .....
9               { "Stud_id":2554, "Stud_name":XYZ, "Stud_marks":25},
10              { "Stud_id":2255, "Stud_name":GHI, "Stud_marks":27}],
11     "status" : "success"
12 },
13 Response Code = 200

```

- 5) **Endpoint (/write, method=POST):** This endpoint writes data entries in a shard in a particular server container. The endpoint expects multiple entries to be written in the server container along with *Shard_id* and the **current index for the shard**. An example request-response pair is shown below.

```

1 Payload Json= {
2     "shard":"sh2",
3     "curr_idx": 507
4     "data": [{ "Stud_id":2255, "Stud_name":GHI, "Stud_marks":27}, ...] /* 5 entries */
5 }
6 Response Json = {
7     "message": "Data entries added",
8     "current_idx": 512, /* 5 entries added */
9     "status" : "success"
10 },
11 Response Code = 200

```

- 6) **Endpoint (/update, method=PUT):** This endpoint updates a particular data entry in a shard in a particular server container. The endpoint expects only one entry to be updated in the server container along with *Shard_id*. An example request-response pair is shown below.

```

1 Payload Json= {
2     "shard":"sh2",
3     "Stud_id":2255,
4     "data": { "Stud_id":2255, "Stud_name":GHI, "Stud_marks":28} /* see marks got updated */
5 }
6 Response Json = {
7     "message": "Data entry for Stud_id:2255 updated",
8     "status" : "success"
9 },
10 Response Code = 200

```

- 7) **Endpoint (/del, method=DELETE):** This endpoint deletes a particular data entry (based on *Stud_id*) in a shard in a particular server container. The endpoint expects only one entry to be deleted in the server container along with *Shard_id*. An example request-response pair is shown below.

```

1 Payload Json= {
2     "shard":"sh1",
3     "Stud_id":2255
4 }
5 Response Json={
6     "message": "Data entry with Stud_id:2255 removed",
7     "status" : "success"
8 },
9 Response Code = 200

```

Finally, write a Dockerfile to containerize the server as an image and make it deployable for the subsequent tasks. Note that two containers can communicate via hostnames in a docker network. Docker provides a built-in DNS service that allows containers to resolve hostnames to the correct IP addresses within the same Docker network.

III. TASK2: IMPROVE THE LOAD BALANCER

You have implemented the load balancer with consistent hashing in the previous assignment. All the previous assumptions and implementation details still hold. However, you must modify and integrate a few features to the load balancer to make it work for assignment 2. In the system diagram shown in Fig. 1, you can observe that the load balancer manages *Stud_id* → *Shard_id* → *Server_id* mapping with two data tables in the **metadata**. The table schemas are as follows:

- 1) *ShardT* (*Stud_id_low*: Number, *Shard_id*: Number, *Shard_size*:Number, *valid_idx*:Number)
- 2) *MapT* (*Shard_id*: Number, *Server_id*: Number)

A. Assumptions

To connect the distributed database assignment with the earlier load balancer assignment, you should realize that you need to maintain consistent hashmaps for each of the shards. The consistent hashmaps can be identified with the *Shard_id*, and the hashmap entries will be populated by the replicas of the shard.

- There are only 4 (sh1, sh2, sh3, sh4) shards in the database.
- Each shard has 3 replicas across the servers.
- There are 6 servers having shards in configuration:

```

1 {
2     "N":6
3     "schema":{"columns":["Stud_id","Stud_name","Stud_marks"],
4         "dtypes":["Number","String","String"]}
5     "shards":[{"Stud_id_low":0, "Shard_id": "sh1", "Shard_size":4096},
6         {"Stud_id_low":4096, "Shard_id": "sh2", "Shard_size":4096},
7         {"Stud_id_low":8192, "Shard_id": "sh3", "Shard_size":4096},
8         {"Stud_id_low":12288, "Shard_id": "sh4", "Shard_size":4096}]
9     "servers":{"Server0":["sh1","sh2"],
10         "Server1":["sh3","sh4"],
11         "Server3":["sh1","sh3"],
12         "Server4":["sh4","sh2"],
13         "Server5":["sh1","sh4"],
14         "Server6":["sh3","sh2"]}
15 }

```

- Thus consistent hashmap for shard sh1 will contain entries for servers {Server0:sh1, Server3:sh1, Server5:sh1}. Similarly each shard would have corresponding hashmap. Other constants for hashmap are as follows:
 - 1) Total number of slots in the consistent hash map ($\#slots$) = 512
 - 2) Number of virtual servers for each server container (K) = $\log(512) = 9$
 - 3) Function for request mapping $H(i)$, & virtual server mapping $\Phi(i, j)$ is what you found works best for your load balancer implementation.

B. Endpoints

The read requests to a particular shard (i) will be load balanced with the consistent hash among all shard (i) replicas across all the server containers. Therefore, parallel read requests can be serviced. In case of write requests, the load balancer ensures consistent write to all the appropriate shard replicas with a mutex lock. Hence, writes are blocking for a particular shard and can only be parallelized across different shard instances (e.g., not replicas). The endpoints of the load balancer are as follows.

- 1) **Endpoint (/init, method=POST):** This endpoint initializes the distributed database across different shards and replicas in the server containers. As administrators, we should provide the configuration of shards and their placements in the set of server containers (place randomly if not specified). An example request-response pair is shown below.

```

1 Payload Json= {
2     "N":3
3     "schema":{"columns":["Stud_id","Stud_name","Stud_marks"],
4         "dtypes":["Number","String","String"]}
5     "shards":[{"Stud_id_low":0, "Shard_id": "sh1", "Shard_size":4096},
6         {"Stud_id_low":4096, "Shard_id": "sh2", "Shard_size":4096},
7         {"Stud_id_low":8192, "Shard_id": "sh3", "Shard_size":4096},]
8     "servers":{"Server0":["sh1","sh2"],
9         "Server1":["sh2","sh3"],

```

```

10         "Server2":["sh1","sh3"]}
11     }
12     Response Json={
13         "message" : "Configured Database",
14         "status" : "success"
15     },
16     Response Code = 200

```

- 2) **Endpoint (/status, method=GET):** This endpoint sends the database configurations upon request. The configuration is set via the `/init` endpoint; thus, a valid response according to the above `/init` example is shown below.

```

1 Response Json= {
2     "N":3
3     "schema":{"columns":["Stud_id","Stud_name","Stud_marks"],
4               "dtypes":["Number","String","String"]}
5     "shards":[{"Stud_id_low":0, "Shard_id": "sh1", "Shard_size":4096},
6               {"Stud_id_low":4096, "Shard_id": "sh2", "Shard_size":4096},
7               {"Stud_id_low":8192, "Shard_id": "sh3", "Shard_size":4096},]
8     "servers":{"Server0":["sh1","sh2"],
9               "Server1":["sh2","sh3"],
10              "Server2":["sh1","sh3"]}
11 }
12 Response Code = 200

```

- 3) **Endpoint (/add, method=POST):** This endpoint adds new server instances in the load balancer to scale up with increasing client numbers in the system. The endpoint expects a JSON payload that mentions the number of new instances, their server names, and the shard placements. An example request and response is below.

```

1 Payload Json= {
2     "n" : 2,
3     new_shards:[{"Stud_id_low":12288, "Shard_id": "sh5", "Shard_size":4096}]
4     "servers" : {"Server4":["sh3","sh5"], /* new shards must be defined */
5                 "Server[5]":["sh2","sh5"]}
6 }
7 Response Json ={
8     "N":5,
9     "message" : "Add Server:4 and Server:58127", /* server id is randomly set in case of
10    Server[5] */
11     "status" : "successful"
12 },
13 Response Code = 200

```

Perform simple sanity checks on the request payload and ensure that servers mentioned in the Payload are less than or equal to newly added instances unless through an error. Server IDs are preferably set randomly to make the consistent hashing algorithm work efficiently. Therefore, you can also define a template for servernames (shown above) where ids are randomly set by the load balancer.

```

1 Payload Json= {
2     "n" : 3, /* wrong input n > len(servers) */
3     new_shards:[{"Stud_id_low":12288, "Shard_id": "sh5", "Shard_size":4096}]
4     "servers" : {"Server4":["sh3","sh5"],
5                 "Server[5]":["sh2","sh5"]}
6 }
7 Response Json ={
8     "message" : "<Error> Number of new servers (n) is greater than newly added instances",
9     "status" : "failure"
10 },
11 Response Code = 400

```

- 4) **Endpoint (/rm, method=DELETE):** This endpoint removes server instances in the load balancer to scale down with decreasing client or system maintenance. The endpoint expects a JSON payload that mentions the number of instances to be removed and their preferred server names in a list. An example request and response is below.

```

1 Payload Json= {
2     "n" : 2,
3     "servers" : ["Server4"]
4 }
5 Response Json ={
6     "message" : {
7         "N" : 3,
8         "servers" : ["Server1", "Server4"] /*See "Server1" is choosen randomly to be
9         deleted along with mentioned "Server4" */
10     },
11     "status" : "successful"
12 },
13 Response Code = 200

```

Perform simple sanity checks on the request payload and ensure that server names mentioned in the Payload are less than or equal to the number of instances to be removed. Note that the server names are preferably mentioned with the delete request. **One can never set the server names. In that case, servers are randomly selected for removal.** However, sending a server list with a greater length than the number of removable instances will result in an error.

```

1 Payload Json= {
2     "n" : 2,          /* wrong input n < len(servers) */
3     "servers" : ["Server1", "Server4", "Server2"]
4 }
5 Response Json = {
6     "message" : "<Error> Length of server list is more than removable instances",
7     "status" : "failure"
8 },
9 Response Code = 400

```

- 5) **Endpoint (/read, method=POST):** Based on the consistent hashing algorithm, this endpoint reads data entries from the shard replicas across all server containers. The endpoint expects a range of *Stud_ids* (low, high) to be read from the distributed database. The load balancer derives all the required shards to be read from the provided *Stud_id* range and accumulates all read results before sending them to the client.

Hint: Due to our indexing (*Stud_id_low*, *Stud_id_low+Shard_size*) rule for each shard, we can easily find the shard entry from the *Stud_id*. (Calculation: *Stud_id_low*+ *Stud_id* % *Shard_size* entry will have the corresponding *Shard_id*).

```

1 Payload Json= {
2     "Stud_id": {"low":1000, "high":8889}
3 }
4 Response Json = {
5     "shards_queried": ["sh1", "sh2", "sh3"]
6     "data" : [{"Stud_id":1000, "Stud_name":PQR, "Stud_marks":23},
7               {"Stud_id":1001, "Stud_name":STV, "Stud_marks":22},
8               ....,
9               {"Stud_id":8888, "Stud_name":ZQN, "Stud_marks":65},
10              {"Stud_id":8889, "Stud_name":BYHS, "Stud_marks":76}],
11     "status" : "success"
12 },
13 Response Code = 200

```

- 6) **Endpoint (/write, method=POST):** This endpoint writes data entries in the distributed database. The endpoint expects multiple entries that are to be written in the server containers. The load balancer schedules each write to its corresponding shard replicas and ensures data consistency using mutex locks for a particular shard and its replicas. The general write work flow will be like: (1) Get *Shard_ids* from *Stud_ids* and group writes for each shard → For each shard Do: (2a)Take mutex lock for *Shard_id* (m) → (2b) Get all servers (set S) having replicas of *Shard_id* (m) → (2c) Write entries in all servers (set S) in *Shard_id* (m) → (2d) Update the *valid_idx* of *Shard_id* (m) in the metadata if writes are successful → (2e) Release the mutex lock for *Shard_id* (m). An example request-response pair is shown.

```

1 Payload Json= {
2     "data": [{"Stud_id":2255, "Stud_name":GHI, "Stud_marks":27},
3             {"Stud_id":3524, "Stud_name":JKBFSFS, "Stud_marks":56}
4             ...
5             {"Stud_id":1005, "Stud_name":YUBAAD, "Stud_marks":100}] /* 100 entries */
6 }
7 Response Json = {
8     "message": "100 Data entries added",
9     "status" : "success"
10 },
11 Response Code = 200

```

- 7) **Endpoint (/update, method=PUT):** This endpoint updates a particular data entry (based on *Stud_id*) in the distributed database. The load balancer retrieves all the shard replicas and their corresponding server instances where the entry has to be updated. The endpoint expects only one entry to be updated. An example request-response pair is shown. *Hint:* Note that the updates are blocking for a particular shard and its replicas to maintain data consistency.

```

1 Payload Json= {
2     "Stud_id": "2255",
3     "data": {"Stud_id":2255, "Stud_name":GHI, "Stud_marks":30} /* see marks got updated */
4 }
5 Response Json = {
6     "message": "Data entry for Stud_id: 2255 updated",
7     "status" : "success"
8 },
9 Response Code = 200

```

- 8) **Endpoint (/del, method=DELETE):** This endpoint deletes a particular data entry (based on *Stud_id*) in the distributed database. The load balancer retrieves all the shard replicas and their corresponding server instances where the entry has to be deleted. The endpoint expects only one entry that is to be deleted. An example request-response pair is shown. *Hint:* Note that the deletes are blocking for a particular shard and its replicas to maintain data consistency.

```

1 Payload Json= {
2     "Stud_id":2255
3 }
4 Response Json = {
5     "message": "Data entry with Stud_id:2255 removed from all replicas",
6     "status" : "success"
7 },
8 Response Code = 200

```

IV. TASK3: ANALYSIS

To analyze the performance of the developed distributed database. You need to perform four subtasks. As mentioned earlier, the design provides scaling in two ways: (i) Read speed with more shard replicas and (ii) Size with more shards and servers.

- A-1 Report the read and write speed for 10000 writes and 10000 reads in the default configuration given in task 2.
- A-2 Increase the number of shard replicas (to 7) from the configuration (init endpoint). Report the write speed down for 10000 writes and read speed up for 10000 reads.
- A-3 Increase the number of Servers (to 10) by adding new servers and increase the number of shards (shard to 6, shard replicas to 8). Define the (init endpoint) configurations according to your choice. Report the write speed up for 10000 writes and read speed up for 10000 reads.
- A-4 Finally, check all the endpoints and ensure their correctness. Manually drop a server container and show that the load balancer spawns a new container and copies the shard entries from other replicas.

APPENDIX

We encourage you to use MYSQL:8.0 to implement the data tables in this assignment. However, you can use any database you want. To run the database and a flask app in the same container, you can follow the method below.

```

1 #Dockerfile
2 #-----
3 FROM mysql:8.0-debian
4
5 COPY deploy.sh /always-initdb.d/          #here the flask app deploy script is copied
6 COPY . /bkr
7 WORKDIR /bkr
8
9 RUN apt-get update
10 RUN apt-get install -y python3
11 RUN apt-get install -y python3-pip
12
13 RUN pip install --upgrade pip
14 RUN pip install -r requirements.txt
15
16 ENV MYSQL_ROOT_PASSWORD="abc"    #host='localhost', user='root',password='abc'
17
18 EXPOSE 5000

```

Whatever is copied in the '/always-initdb.d/' folder will be run after the database is initialized. Therefore, the Flask app can connect to the local database with root credentials (`mysql.connector.connect(host='localhost',user='root',password='abc')`).

```

1 #deploy.sh file
2 #-----
3 #!/bin/bash
4 python3 flaskApp.py &

```

For reference, you can observe how the Dockerfile and deploy.sh file is coded for <https://github.com/prasenj52282/shardQ/tree/main/broker>. For other databases (i.e., MongoDB, Cassandra, etc.), similar always-init scripts can be run when the database initializes in every startup. Feel free to choose any database of your choice.

A. Grading Scheme

- TASK1: Server - 40 %
- TASK2: Improve the load balancer - 30 %
- TASK3: Analysis - 20 %
- Documentation & Code Optimizations - 10 %

REFERENCES

- [1] Docker, "What is a container?." <https://docs.docker.com/guides/get-started/>, 2024.
- [2] Docker, "Use docker compose." https://docs.docker.com/get-started/08_using_compose/, 2024.
- [3] makefiletutorial, "Learn makefiles with the tastiest examples." <https://makefiletutorial.com/>, 2024.
- [4] J. Li, Y. Nie, and S. Zhou, "A dynamic load balancing algorithm based on consistent hash," in *2018 2nd IEEE Advanced Information Management, Communicates, Electronic and Automation Control Conference (IMCEC)*, pp. 2387–2391, 2018.
- [5] T. Roughgarden and G. Valiant, "Cs168: The modern algorithmic toolbox lecture 1: Introduction and consistent hashing." <https://web.stanford.edu/class/cs168/l1.pdf>, 2022.
- [6] DockerHub, "MySQL:8.0 docker image source in dockerhub." https://hub.docker.com/_/mysql, 2024.