# ENGR-E516 Assignment – Map-Reduce with VMs and FaaS

## Part 1 – Map Reduce on VMs

### Overview

In Part 1, the Map Reduce implementation from Assignment 2 has been deployed on GCP Compute Engine instances. Each mapper/reducer VM is launched/destroyed only from the master. The master also creates the KV store VM for storage. All communication between VMs is TCP socket-based.

### Steps to Execute/Test Implementation

- The master VM is already created and it's public IP is 34.140.99.15
  There is a Flask server on port 8081 with below 2 endpoints. It is recommended to use an app like Postman for testing these APIs.

    - **To launch map reduce**
      **Note:** This process will take roughly 5-8 minutes to complete as the VMs need to be created and the execution environment/dependencies need to set up for all mapper and reducer VMs.

      Send below POST request at endpoint "/launch_map_reduce"

      cURL request
      curl --location --request POST '34.140.99.15:8081/launch_map_reduce'

    - **To fetch final output**
      Send below GET request at endpoint "/final_output"
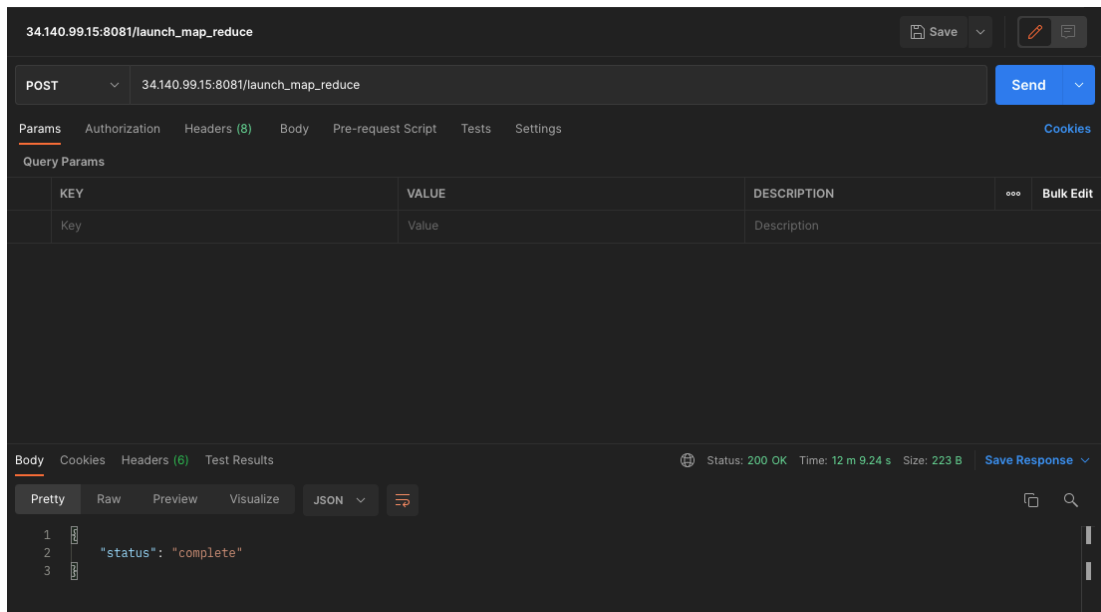
      cURL request
      curl --location --request GET '34.140.99.15:8081/final_output'

# Test Example using Postman

After the entire map-reducer workflow is complete, the status of "complete" will be seen in the response. Then, one can send a GET request on the "/final_output" API to fetch the result.

## Launch Map Reduce

## Fetch Result



34.140.99.15:8081/final_output

GET     34.140.99.15:8081/final_output     Send

Params   Authorization   Headers (7)   Body   Pre-request Script   Tests   Settings     Cookies

Query Params

| KEY | VALUE | DESCRIPTION | | Bulk Edit |
|-----|-------|-------------|---|----------|
| Key | Value | Description | | |

Body   Cookies   Headers (6)   Test Results     Status: 200 OK   Time: 1122 ms   Size: 692.81 KB   Save Response

Pretty   Raw   Preview   Visualize   JSON

```
 1  {
 2      "a": [
 3          "doc2.txt",
 4          "doc3.txt",
 5          "doc1.txt",
 6          "doc4.txt"
 7      ],
 8      "a1": [
 9          "doc1.txt"
10      ],
11      "a3": [
12          "doc1.txt"
13      ],
14      "a4": [
15          "doc1.txt"
16      ],
17      "a5": [
18          "doc1.txt"
19      ],
20      "aa": [
21          "doc4.txt"
22      ],
23      "aaron": [
24          "doc2.txt"
25  ]
```

# Implementation

KV Store, Mapper and Reducer VMs are created, set up and triggered by the master VM.



**Barrier**: Reducer VMs wait for Mapper VMs to complete their task.

```
Pseudo-terminal will not be allocated because stdin is not a terminal.
[MASTER] Connection request accepted from mapper ('10.132.0.14', 36946)
[MASTER] mapper3 task completed

[MASTER] ACK received from all 3 mappers

[MASTER] --------- BARRIER (waiting for mappers to complete) ----------


[MASTER] All 3 mapper tasks are complete...

[MAPPER - mapper3] mapper3 has started...
[MAPPER - mapper3] Inverted index mapping started...
Waiting for operation to finish...
done.
Waiting for operation to finish...
done.
**** PWD !!!! *** is /home/niravraje/gcp-map-reduce
Setting up 2 reducers. Sleeping 30 seconds to allow the VMs to accept SSH on port 22
```

Reducers are launched only after the mappers are complete.



Final output is stored in the KV store itself. The /final_output API is provided for external users to fetch the output.

The cleanup script is run by the master after process is complete to delete all mappers and reducers.

The KV Store is not deleted by the script in order to allow GET operations on final output.

```
[MASTER] ACK received from all 2 reducers
[MASTER] Generating final output file & writing to ./gcp-map-reduce/kv-data-store/final-output...
+ MAPPER_COUNT=3
+ REDUCER_COUNT=2
+ INSTANCE_ZONE=europe-west1-b
+ (( i=1 ))
+ (( i<=3 ))
+ MAPPER_INSTANCE_NAME=mapper1
+ echo 'Deleting mapper1 VM'
Deleting mapper1 VM
+ (( i++ ))
+ (( i<=3 ))
+ MAPPER_INSTANCE_NAME=mapper2
+ echo 'Deleting mapper2 VM'
Deleting mapper2 VM
+ (( i++ ))
+ (( i<=3 ))
+ MAPPER_INSTANCE_NAME=mapper3
+ echo 'Deleting mapper3 VM'
Deleting mapper3 VM
+ (( i++ ))
+ (( i<=3 ))
+ (( i=1 ))
+ (( i<=2 ))
+ REDUCER_INSTANCE_NAME=reducer1
+ echo 'Deleting reducer1 VM'
Deleting reducer1 VM
+ (( i++ ))
+ (( i<=2 ))
+ REDUCER_INSTANCE_NAME=reducer2
+ echo 'Deleting reducer2 VM'
Deleting reducer2 VM
+ (( i++ ))
+ (( i<=2 ))
+ gcloud compute instances delete mapper2 --zone=europe-west1-b --delete-disks=all --quiet
+ gcloud compute instances delete mapper3 --zone=europe-west1-b --delete-disks=all --quiet
+ gcloud compute instances delete reducer1 --zone=europe-west1-b --delete-disks=all --quiet
+ gcloud compute instances delete reducer2 --zone=europe-west1-b --delete-disks=all --quiet
```

# Part 2 – Map Reduce using FaaS (Google Cloud Functions)

## Overview

In this part, the entire map reduce workflow is processed using FaaS.

Below are the cloud functions created:

- Master
- Mapper
- Reducer
- Streaming master trigger – Used for monitoring any uploads to the "raw-dataset" bucket. Any new file upload triggers the master process.
- Web UI – A cloud function serving the inverted index search UI
- UI Handler – Accepts the word entered by the user on the Web UI, searches for the word in the final output and returns the result to display on the UI.



## Implementation

### Web UI
The UI can be accessed at: https://web-ui-cf-gf2co4zquq-uc.a.run.app



Currently there are 4 documents in the raw-dataset bucket.

To search another word, please refresh the page. This is a basic UI to test implementation and can be improved further to indicate "Word not found" or show a progress indicator.



## Streaming Search

On creating a new file or uploading one for the "mr-raw-dataset" bucket, the below CF will trigger master as shown below.



Source code

```
 1   import requests
 2
 3   MASTER_FAAS_URL = "https://master-cf-gf2co4zquq-uc.a.run.app"
 4
 5   def hello_gcs(event, context):
 6       """Triggered by a change to a Cloud Storage bucket.
 7       Args:
 8           event (dict): Event payload.
 9           context (google.cloud.functions.Context): Metadata for the event.
10       """
11       file = event
12       print(f"Processing file: {file['name']}.")
13       requests.get(MASTER_FAAS_URL)
14       return {"status": "master_triggered"}
15
```

As we can see below, when "doc3.txt" and "doc4.txt" were uploaded, the master CF was triggered.

```
>  ⚙   2022-12-05 05:40:51.096 EST    streaming-master-trigger-cf  5n2m4ics1aws   Function execution started
>  ☀   2022-12-05 05:40:53.594 EST    streaming-master-trigger-cf  5n2m4ics1aws   Processing file: doc3.txt.
>  ⚙   2022-12-05 05:41:21.811 EST    streaming-master-trigger-cf  5n2m4ics1aws   Function execution took 30719 ms, finished with status: 'ok'
>  ⚙   2022-12-05 13:48:03.489 EST    streaming-master-trigger-cf  wz0k3h9akt76   Function execution started
∨  ☀   2022-12-05 13:48:06.956 EST    streaming-master-trigger-cf  wz0k3h9akt76   Processing file: doc4.txt.
   ▾  {
        insertId: "638e3ce6000e99861eb4dc49"
      ▸ labels: {2}
        logName: "projects/nirav-raje-fall2022/logs/cloudfunctions.googleapis.com%2Fcloud-functions"
        receiveTimestamp: "2022-12-05T18:48:06.964503871Z"
      ▸ resource: {2}
        textPayload: "Processing file: doc4.txt."
        timestamp: "2022-12-05T18:48:06.956806Z"
        trace: "projects/nirav-raje-fall2022/traces/92c42e21ad80e71a32bedf5f39124b4d"
      }
>  ⚙   2022-12-05 13:48:45.912 EST    streaming-master-trigger-cf  wz0k3h9akt76   Function execution took 42429 ms, finished with status: 'ok'
ⓘ   No newer entries found matching current filter.
```

## Master Function – Orchestration & Barrier Implementation

The mapper and reducer cloud functions are simply called in parallel from the master.
Parallel invocation has been achieved using multithreading in Python.
Barrier has been implemented using thread.join() in a loop.

```
MAPPER_FAAS_URL = "https://mapper-cf-gf2co4zquq-uc.a.run.app"
REDUCER_FAAS_URL = "https://reducer-cf-gf2co4zquq-uc.a.run.app"

def mapper_trigger(mapper_id):
    response = requests.get(MAPPER_FAAS_URL, params={"mapper_id": str(mapper_id)})

def reducer_trigger(reducer_id):
    response = requests.get(REDUCER_FAAS_URL, params={"reducer_id": str(reducer_id)})
```

```python
mapper_threads = []
for i in range(1, MAPPER_COUNT+1):
    mapper_id = f"mapper{i}"
    curr_thread = Thread(target=mapper_trigger, args=(mapper_id,))
    mapper_threads.append(curr_thread)
    curr_thread.start()

# Barrier: wait for all mapper cloud functions to complete
for curr_thread in mapper_threads:
    curr_thread.join()

reducer_threads = []
for i in range(1, REDUCER_COUNT+1):
    reducer_id = f"reducer{i}"
    curr_thread = Thread(target=reducer_trigger, args=(reducer_id,))
    reducer_threads.append(curr_thread)
    curr_thread.start()

# Wait for all reducer cloud functions to complete
for curr_thread in reducer_threads:
    curr_thread.join()

# Combine all reducer output files into a single final output file
final_output = combine_reducer_output_files()
upload_final_output_to_cloud_storage(final_output)
```

The below script snippets can be used to create the cloud functions from the gcloud CLI instead of Console. Note: Each snippet needs to be run from the appropriate folder (please see attached codebase for reference)

```
# Deploy Master
gcloud functions deploy master-cf \
--gen2 \
--runtime=python310 \
--region=us-central1 \
--source=. \
--entry-point=master_init \
--trigger-http \
--allow-unauthenticated

# Deploy mapper faas
gcloud functions deploy mapper-cf \
--gen2 \
--runtime=python310 \
--region=us-central1 \
--source=. \
--entry-point=mapper_init \
--trigger-http \
--allow-unauthenticated

# Deploy reducer faas
gcloud functions deploy reducer-cf \
--gen2 \
--runtime=python310 \
--region=us-central1 \
--source=. \
--entry-point=reducer_init \
--trigger-http \
--allow-unauthenticated

# Deploy UI handler
gcloud functions deploy ui-handler-cf \
--gen2 \
--runtime=python310 \
--region=us-central1 \
--source=. \
--entry-point=handle_request \
--trigger-http \
--allow-unauthenticated

# Deploy UI
gcloud functions deploy web-ui-cf \
--gen2 \
--runtime=python310 \
--region=us-central1 \
--source=. \
--entry-point=launch_ui \
--trigger-http \
--allow-unauthenticated
```