# Big Data Lab 6 - Assignment

## Vasudev Gupta (ME18B182)

1) Following workflow is followed for the setting up Pub/Sub:
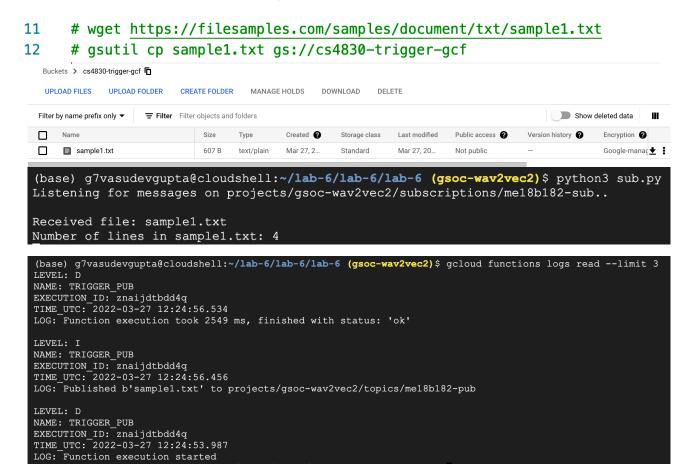
A topic is created and the subscription is setup using following commands.

```
1    # gcloud pubsub topics create me18b182-pub
2    # gcloud pubsub subscriptions create me18b182-sub --topic me18b182-pub
```

Then, a google cloud function is deployed using the following command:



Then, script **sub.py** is run and following command is run to download the dataset and copying it to the GCS bucket. After few seconds, we got the response as shown in third image below.

```
11   # wget https://filesamples.com/samples/document/txt/sample1.txt
12   # gsutil cp sample1.txt gs://cs4830-trigger-gcf
```

Following code is used to trigger the cloud function upon adding file to bucket:

```python
16    def TRIGGER_PUB(data, context):
17        from google.cloud import pubsub_v1
18
19        project_id = "gsoc-wav2vec2"
20        topic_id = "me18b182-pub"
21
22        # Create a publisher client to topic
23        publisher = pubsub_v1.PublisherClient()
24        topic_path = publisher.topic_path(project_id, topic_id)
25
26        # Publish path to file in the topic & Data must be a bytestring
27        pub_data = data['name'].encode("utf-8")
28
29        # Client returns future when publisher publish.
30        future = publisher.publish(topic_path, pub_data)
31
32        print(f"Published {pub_data} to {topic_path}")
```

Following is the code contained in **sub.py** :

```python
6     project_id = "gsoc-wav2vec2"
7     subscription_id = "me18b182-sub"
8
9     # Create a subscriber client
10    subscriber = pubsub_v1.SubscriberClient()
11    subscription_path = subscriber.subscription_path(project_id, subscription_id)
12
13    def callback(message):
14        # Decode the filename back to string
15        filename = message.data.decode("utf-8")
16        print(f"Received file: {filename}")
17
18        # Create a google storage client to read the uploaded file from bucket.
19        storage_client = storage.Client()
20        bucket_name = "cs4830-trigger-gcf"
21        bucket = storage_client.get_bucket(bucket_name)
22
23        blob = bucket.blob(filename).download_as_string().decode('utf-8')
24
25        # Splits the text based on \n
26        lines = blob.split('\n')
27        print(f"Number of lines in {filename}: {len(lines)}")
28
29        message.ack()
30
31    # Default subscriber is a pull subscriber
32    pull_sub_future = subscriber.subscribe(subscription_path, callback=callback)
33    print(f"Listening for messages on {subscription_path}..\n")
34
35    try:
36        # The subscriber listens indefinitely
37        pull_sub_future.result()
38    except KeyboardInterrupt:
39        pull_sub_future.cancel()
```

2) A subscription can use either the pull or push mechanism for message delivery.

**Pull Subscription** - In this system, the subscriber application initiates requests to the Pub/Sub server to retrieve messages.
1. The pull method is called explicitly by the subscribing application, which requests that messages be sent.
2. The message (or an error if the queue is empty) and an ack ID are returned by the Pub/Sub server.
3. The subscriber expressly invokes the acknowledge method and acknowledges receipt using the returned ack ID.

**Push Subscription -** In push delivery, Pub/Sub initiates requests to your subscriber application to deliver messages.
1. Each message is sent to the subscriber application as an HTTPS request to a pre-configured endpoint by the Pub/Sub server.
2. The message is acknowledged by the endpoint, which returns an HTTP success code. A message should be resent if it receives a non-success answer.

Scenarios where pull subscription is favourable -
- When there is a large volume of messages.
- The efficiency and throughput of message processing are critical.
- Public HTTPS endpoint, with a non-self-signed SSL certificate, is not feasible to setup

Scenarios where push subscription is favourable -
- Multiple topics must be processed by the same webhook.
- App Engine Standard and Cloud Functions subscribers.
- Environments where Google Cloud dependencies (such as credentials and the client library) are not feasible to set up.