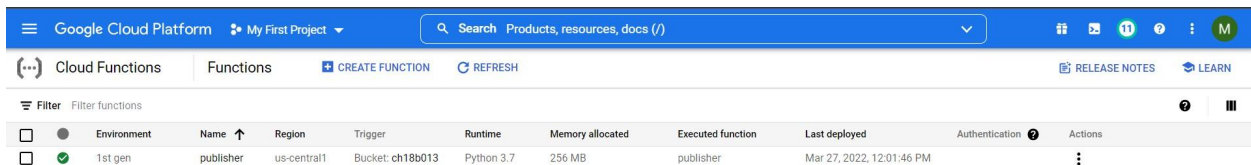**CS4830 - BIG DATA LABORATORY**
**LAB - 6**
**VISHAL RISHI MK - CH18B013**

**Q1. In this assignment, you will count the number of lines in a file uploaded to the GCS bucket in real-time by using Google Cloud Functions and Pub/Sub.**

**a)** The file *'sample1.txt'* is downloaded and saved in the terminal.
**b)** The Google Cloud Function '*publisher'* gets triggered whenever a file is added to a bucket and publishes the file name to a topic in Pub/Sub. It is given in the '*main.py'* file.
**c)** *'subscriber.py'* acts as a subscriber to this topic and prints out the number of lines in the file in real-time.



*Figure 1: Deploying the GCF*



*Figure 2: Deploying the GCF*



*Figure 3: Output of the subscriber. We count the number of lines in real-time in the sample1.txt file using this code. Hence, we have created the subscriber to the topic and we have calculated the number of lines in the sample.txt document.*

**Q2. There are two kinds of subscribers - pull and push subscribers. What are the differences between the two and when would you prefer one over the other?**

Subscriptions are named resources representing the stream of messages from a single, specific topic, to be delivered to the subscribing application. A publisher application creates and sends messages to a topic. Subscriber applications create a subscription to a topic to receive messages from it. There are two types of subscriptions that one can choose from: pull subscriptions and push subscriptions.

**Push Subscription:**
With a push subscription, the Publisher propagates changes to a Subscriber without a request from the Subscriber. Changes can be pushed to Subscribers on demand, continuously, or on a scheduled basis. Push Subscription is used when:
• Data needs to be synchronized continuously or on a frequently recurring schedule.
• Publications require near real-time movement of data.
• Most often used with snapshot and transactional replication.

**Pull Subscription:**
With a pull subscription, the Subscriber requests changes made at the Publisher. Pull subscriptions allow the user at the Subscriber to determine when the data changes are synchronized. Pull Subscription is used when:
• Data needs to be synchronized on-demand or on a schedule rather than continuously.
• The publication has a large number of Subscribers, and/or it would be too resource-intensive to run all the agents at the Distributor.
• Subscribers are autonomous, disconnected, and/or mobile. Subscribers will determine when they will connect and synchronize changes.
• Most often used with merge replication.