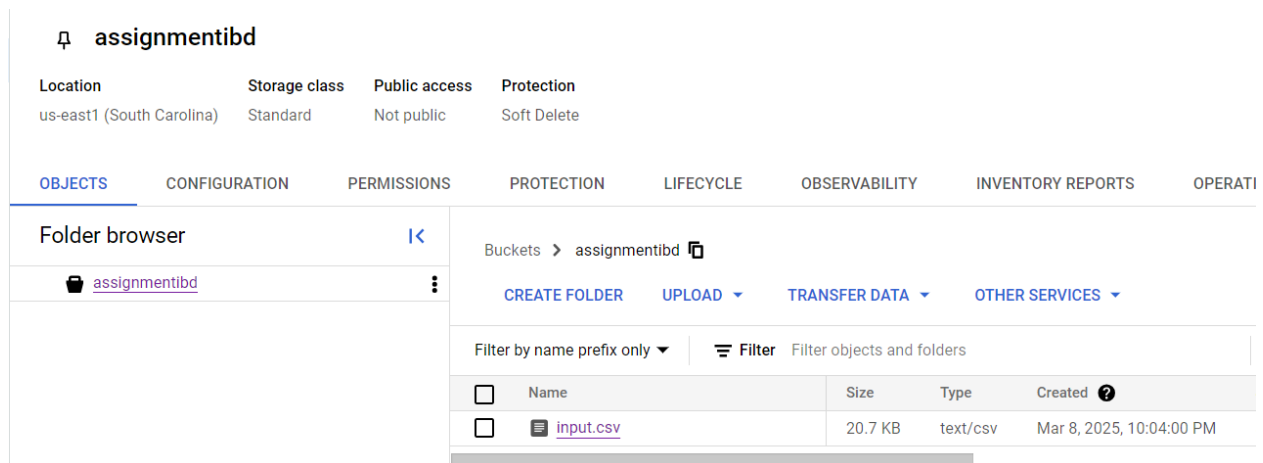# Question: Write a Producer that reads from that file, breaks the data into, and writes to Kafka and a Spark Streaming consumer that reads from the topic and emits count of rows.

Submitted by Shramana Sinha, 23f1002703

## 1. Implementation Steps

### 1.1 Google Cloud Storage Configuration

- Created a bucket named `assignmentibd`
- Uploaded an input file `input.csv` containing customer data with 1,000 records



*Screenshot: The gcs bucket with the input.csv file*

### 1.2 Network Configuration

- Created a firewall rule to allow Kafka traffic on port 9092:

```
Unset
gcloud compute firewall-rules create allow-kafka --allow tcp:9092
--description="Allow Kafka traffic" --direction=INGRESS
```

### 1.3 Kafka Server Setup

- Launched a VM instance with Ubuntu 20.04 LTS
- Installed Java and Python dependencies:

```
Unset
sudo apt update && sudo apt install -y openjdk-8-jdk wget python3-pip
pip install kafka-python google-cloud-storage
```

- Downloaded and extracted Apache Kafka 3.9.0:

```
Unset
wget https://dlcdn.apache.org/kafka/3.9.0/kafka_2.13-3.9.0.tgz
tar -xzf kafka_2.13-3.9.0.tgz
```

- Configured Kafka server to accept external connections:

```
Unset
# In kafka_2.13-3.9.0/config/server.properties
listeners=PLAINTEXT://0.0.0.0:9092
advertised.listeners=PLAINTEXT://<VM_EXTERNAL_IP>:9092
```

- Started Zookeeper and Kafka server as daemon processes:

```
Unset
kafka_2.13-3.9.0/bin/zookeeper-server-start.sh -daemon config/zookeeper.properties
kafka_2.13-3.9.0/bin/kafka-server-start.sh -daemon config/server.properties
```

- Uploaded the `producer.py` script and started its execution:

```
Unset
python3 producer.py
```

## VM instances

Filter    Enter property name or value

| | Status | Name ↑ | Zone | Recommendations | In use by | Internal IP | External IP | Connect | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | ✅ | cluster-aa48-m | us-central1-c | | | 10.128.15.205 (nic0) | 35.192.209.11 (nic0) | SSH ▾ | ⋮ |
| ☐ | ✅ | cluster-aa48-w-0 | us-central1-c | | | 10.128.15.204 (nic0) | 35.194.0.122 (nic0) | SSH ▾ | ⋮ |
| ☐ | ✅ | cluster-aa48-w-1 | us-central1-c | | | 10.128.15.203 (nic0) | 35.223.229.250 (nic0) | SSH ▾ | ⋮ |
| ☐ | ✅ | instance-20250309-042530 | us-central1-c | | | 10.128.15.193 (nic0) | 34.58.97.69 (nic0) | SSH ▾ | ⋮ |

*Screenshot: The vm (instance-20250309-042530) for running the kafka and producer.py*

## 1.4 Dataproc Cluster Setup

- Created a Dataproc cluster with external IP access enabled
- Uploaded the `consumer.py` script.
- Submitted the spark job

```Unset
spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.4
consumer.py
```

| | |
|---|---|
| **Name** | cluster-aa48 |
| **Cluster UUID** | 60d048e7-e81c-4ecf-b4f5-477097637dd1 |
| **Type** | Dataproc Cluster |
| **Status** | ✅ Running |

| MONITORING | JOBS | VM INSTANCES | CONFIGURATION | WEB INTERFACES |
|---|---|---|---|---|

≡ Filter   Filter instances

| | Name | Role |
|---|---|---|
| ✅ | cluster-aa48-m | Master |
| ✅ | cluster-aa48-w-0 | Worker |
| ✅ | cluster-aa48-w-1 | Worker |

*Screenshot: The dataproc cluster for running the consumer.py*

## 2. Data Producer Implementation

The producer application performs the following tasks:

- Downloads the CSV file from Google Cloud Storage

```Python
def download_from_gcs(bucket_name, source_blob_name, destination_file_name):
    storage_client = storage.Client()
    bucket = storage_client.bucket(bucket_name)
    blob = bucket.blob(source_blob_name)
    blob.download_to_filename(destination_file_name)
    print(
        f"Downloaded gs://{bucket_name}/{source_blob_name} to
{destination_file_name}"
    )
```

- Reads and parses the customer data

```python
def read_local_file(file_path):
    rows = []
    with open(file_path, "r") as file:
        csv_reader = csv.DictReader(file, fieldnames=["customer_id", "name",
"city"])
        for row in csv_reader:
            rows.append(row)
    return rows
```

- Streams the data to Kafka in batches of 10 records

```python
records_to_send = len(batch)
batch = batch[:records_to_send]
print(f"Sending batch {i+1} with {len(batch)} records to Kafka...")
for record in batch:
        producer.send(KAFKA_TOPIC, value=record)
        sent_count += 1
producer.flush()
```

- Implements a controlled flow with 10-second pauses between batches

```python
if i < len(batches) - 1:
        print(f"Sleeping for {sleep_seconds} seconds...")
        time.sleep(sleep_seconds)
```

*Screenshot: The producer sending batches of 10 records to kafka every 10 seconds*

# 3. Data Consumer Implementation

The Spark Streaming consumer application performs the following functions:

- Establishes a connection to the Kafka topic and creates a streaming DataFrame from the Kafka source

```python
df = (
        spark.readStream.format("kafka")
        .option("kafka.bootstrap.servers", KAFKA_BOOTSTRAP_SERVERS)
        .option("subscribe", KAFKA_TOPIC)
        .option("startingOffsets", "earliest")
        .load()
    )
```

- Processes the incoming data by casting the binary values to string and getting the timestamp from the message.

```python
value_df = df.selectExpr("CAST(value AS STRING)", "timestamp")
```

- Counts records in 10-second windows with 5-second sliding intervals

```Python
windowedCounts = value_df.groupBy(
        window(col("timestamp"), "10 seconds", "5 seconds")
    ).count()
```

- Outputs results to the console in update mode

```Python
query = (
        windowedCounts.writeStream.outputMode("update")
        .format("console")
        .option("truncate", "false")
        .start()
    )
```



*Screenshot: The output of the consumer.py showing the count of records for 10 seconds windows which are sliding every 5 seconds*