**Name: Omkar Rajesh Desai**
**Email: omkarbkishor@gmail.com**

# Real-Time Advertisement Data Aggregation Documentation

**Tools Used:**
1. Python3 – PyCharm Community Edition
2. Docker
3. Kafka
4. Zookeeper
5. Spark
6. Jupyter
7. Cassandra

**Files Attached:**
1. Documentation.pdf – This file
2. producer.py – file for producing ads_data to kafka topic(i.e ads_data).
3. Docker-compose.yml – For setting up the environment for Kafka, zookeeper, spark, Cassandra.
4. streaming_code.py – which contain streaming code.
5. schema.avsc – which contain Avro-schema used.

providing **GitHub** link for more information

**Process and File Descriptions:**

**Step 1:**

I first created a docker-compose file for setting up the environment for project by considering all the required services for the project and make sure all the required ports for each service are exposed and make sure to keep all the containers connected to the same network. I also created the service for control center and schema registry control center which will provide web UI for the schema registry for our Avro schema which is depend on schema registry. Schema registry help us to register our avro_schema.

```yaml
      KAFKA_TRANSACTION_STATE_LOG_MIN_ISR: 1
      KAFKA_TRANSACTION_STATE_LOG_REPLICATION_FACTOR: 1
      KAFKA_JMX_PORT: 9101
      KAFKA_JMX_HOSTNAME: localhost
      KAFKA_CONFLUENT_SCHEMA_REGISTRY_URL: http://schema-registry:8081
      CONFLUENT_METRICS_REPORTER_BOOTSTRAP_SERVERS: broker:29092
      CONFLUENT_METRICS_REPORTER_TOPIC_REPLICAS: 1
      CONFLUENT_METRICS_ENABLE: 'false'
      CONFLUENT_SUPPORT_CUSTOMER_ID: 'anonymous'
    networks:
      confluent:
        aliases:
          - broker
    healthcheck:
      test: ['CMD', 'bash', '-c', 'nc -z localhost 9092']
      interval: 10s
      timeout: 5s
      retries: 5

  schema-registry:
    image: confluentinc/cp-schema-registry:7.4.0
    hostname: schema-registry
    container_name: schema-registry
    depends_on:
      broker:
        condition: service_healthy
    ports:
      - "8081:8081"
    environment:
      SCHEMA_REGISTRY_HOST_NAME: schema-registry
      SCHEMA_REGISTRY_KAFKASTORE_BOOTSTRAP_SERVERS: 'broker:29092'
      SCHEMA_REGISTRY_LISTENERS: http://0.0.0.0:8081
    networks:
      confluent:
        aliases:
          - schema-registry
    healthcheck:
      test: ['CMD', 'curl', '-f', 'http://localhost:8081']
```

**Fig.: Docker-compose File**

```yaml
version: '3.8'
services:
  zookeeper:
    image: confluentinc/cp-zookeeper:7.4.0
    hostname: zookeeper
    container_name: zookeeper
    ports:
      - "2181:2181"
    environment:
      ZOOKEEPER_CLIENT_PORT: 2181
      ZOOKEEPER_TICK_TIME: 2000
    healthcheck:
      test: ['CMD', 'bash', '-c', "echo 'ruok' | nc localhost 2181"]
      interval: 10s
      timeout: 5s
      retries: 5
    networks:
      confluent:
        aliases:
          - zookeeper

  broker:
    image: confluentinc/cp-server:7.4.0
    hostname: broker
    container_name: broker
    depends_on:
      zookeeper:
        condition: service_healthy
    ports:
      - "9092:9092"
      - "9101:9101"
    environment:
      KAFKA_BROKER_ID: 1
      KAFKA_ZOOKEEPER_CONNECT: 'zookeeper:2181'
      KAFKA_LISTENER_SECURITY_PROTOCOL_MAP: PLAINTEXT:PLAINTEXT,PLAINTEXT_HOST:PLAINTEXT
      KAFKA_ADVERTISED_LISTENERS: PLAINTEXT://broker:29092,PLAINTEXT_HOST://localhost:9092
      KAFKA_METRIC_REPORTERS: io.confluent.metrics.reporter.ConfluentMetricsReporter
      KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
      KAFKA_GROUP_INITIAL_REBALANCE_DELAY_MS: 0
      KAFKA_CONFLUENT_LICENSE_TOPIC_REPLICATION_FACTOR: 1
      KAFKA_CONFLUENT_BALANCER_TOPIC_REPLICATION_FACTOR: 1
```

**Fig.: Docker-compose File**

```yaml
    control-center:
      image: confluentinc/cp-enterprise-control-center:7.4.0
      hostname: control-center
      container_name: control-center
      depends_on:
        broker:
          condition: service_healthy
        schema-registry:
          condition: service_healthy
      ports:
        - "9021:9021"
      environment:
        CONTROL_CENTER_BOOTSTRAP_SERVERS: 'broker:29092'
        CONTROL_CENTER_SCHEMA_REGISTRY_URL: "http://schema-registry:8081"
        CONTROL_CENTER_REPLICATION_FACTOR: 1
        CONTROL_CENTER_INTERNAL_TOPICS_PARTITIONS: 1
        CONTROL_CENTER_MONITORING_INTERCEPTOR_TOPIC_PARTITIONS: 1
        CONFLUENT_METRICS_TOPIC_REPLICATION: 1
        CONFLUENT_METRICS_ENABLE: 'false'
        PORT: 9021
      networks:
        confluent:
          aliases:
            - control-center

  spark-master:
    image: bitnami/spark:latest
    command: bin/spark-class org.apache.spark.deploy.master.Master
    ports:
      - "9090:8080"
      - "7077:7077"
    networks:
      confluent:
        aliases:
          - spark-master
```

**Fig.: Docker-compose File**

```yaml
  spark-worker:
    image: bitnami/spark:latest
    command: bin/spark-class org.apache.spark.deploy.worker.Worker spark://spark-master:7077
    depends_on:
      - spark-master
    environment:
      SPARK_MODE: worker
      SPARK_WORKER_CORES: 2
      SPARK_WORKER_MEMORY: 1g
      SPARK_MASTER_URL: spark://spark-master:7077
    ports:
      - "8085:8081"
    networks:
      confluent:
        aliases:
          - spark-worker

  ed-pyspark-jupyter:
    image: jupyter/pyspark-notebook:latest
    user: root
    container_name: ed-pyspark-jupyter-lab
    ports:
      - "8888:8888"
      - "4040:4040"
    environment:
      JUPYTER_ENABLE_LAB: "yes"
      JUPYTER_PORT: 8888
      SPARK_UI_PORT: 4040
      GRANT_SUDO: yes
    volumes:
      - ./notebooks:/home/jovyan/work
    networks:
      confluent:
        aliases:
          - ed-pyspark-jupyter

  cassandra:
    image: cassandra:latest
    container_name: cassandra
```

**Fig.: Docker-compose File**

<u>**Step 2:**</u>

After setting up the containers checked all the contianers are up and running then create the
**ads_data** topic with 2 partitions on kafka using control center web UI after that create Avro
schema file for schema registry then write logic to produce messages on kafka topic.

```
producer.py ×
1   import random
2   import time
3   from datetime import datetime, timedelta
4   from confluent_kafka import Producer
5   from avro import schema
6   import avro.io
7   import io
8
9
10  # Define the Avro schema as a string
11  avro_schema_str = """
12  {
13    "type": "record",
14    "name": "test",
15    "namespace": "spark_stream_test",
16    "fields": [
17      {"name": "ad_id", "type": "string"},
18      {"name": "timestamp", "type": "string"},
19      {"name": "clicks", "type": "int"},
20      {"name": "views", "type": "int"},
21      {"name": "cost", "type": "double"}
22    ]
23  }
24  """
25
26  # Kafka parameters
27  schema_registry_url = "http://localhost:8081"
28  kafka_broker = "localhost:9092"
29  topic = "ads_data"
30
31
32  # Initialize Kafka Producer
33  producer = Producer({'bootstrap.servers': kafka_broker})
34
35  avro_schema = schema.parse(avro_schema_str)
36
```

```
1 usage   ± Omkar Desai
8   def generate_data():
        ad_id = '123' + str(random.randint( a: 1,  b: 51))
        timestamp = (datetime.now() - timedelta(days=random.randint( a: 0,  b: 30))).isoformat()
        clicks = random.randint( a: 0,  b: 100)
        views = random.randint( a: 0,  b: 500)
        cost = round(random.uniform( a: 5,  b: 100), 2)

        record = {
            "ad_id": ad_id,
            "timestamp": timestamp,
            "clicks": clicks,
            "views": views,
            "cost": cost
        }

        return record


    # Produce messages to Kafka
    for _ in range(100):  # Produce 100 messages
        message_writer = avro.io.DatumWriter(avro_schema)
        message_bytes_writer = io.BytesIO()
        message_encoder = avro.io.BinaryEncoder(message_bytes_writer)
        data = generate_data()
        message_writer.write(data, message_encoder)
        message_raw_bytes = message_bytes_writer.getvalue()

        # Serialize the data using AvroSerializer
        producer.produce(topic, message_raw_bytes)
        producer.flush()
        time.sleep(1)
```

**Fig.: Producer Code**

**Step 3:**

I then went ahead and performed the necessary transformations/queries on produce data using spark structured streaming import all the necessary modules by first install them in container using pip command after that mention all the required jar files in config of spark session such as spark-sql-kafka, spark-avro, spark-cassandra then in code done window_based aggregation with window duration of 1 minute and sliding interval of 30 seconds. Then create function for processing each row's data with existing data in cassandra table then used forEachBatch function to apply this function for each batch and for deserializing messages I used from_avro function.

```
streaming_code.py ×
1   from pyspark.sql import SparkSession
2   from pyspark.sql.functions import col, sum, avg, window, current_timestamp
3   from pyspark.sql.avro.functions import from_avro
4   from cassandra.query import SimpleStatement
5   from cassandra.cluster import Cluster, NoHostAvailable
6
7   # Initialize Spark session
8   spark = SparkSession.builder \
9       .appName("KafkaAvroConsumer") \
10      .config("spark.cassandra.connection.host", "cassandra") \
11      .config("spark.cassandra.connection.port", "9042") \
12      .config("spark.jars.packages",
13          "org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.0,org.apache.spark:spark-avro_2.12:3.5.0,com.datastax.spark:spark-cassandra-connector_2.12:3.5.0") \
14      .getOrCreate()
15
16  # Kafka configuration
```

```python
# Kafka configuration
kafka_bootstrap_servers = 'broker:29092'
kafka_topic = 'test'

# Read Avro schema from file
with open('schema.avsc', mode='r') as file:
    schema_avro = file.read()

# Stream data from Kafka
kafka_df = spark.readStream \
    .format("kafka") \
    .option(key: "kafka.bootstrap.servers", kafka_bootstrap_servers) \
    .option(key: "subscribe", kafka_topic) \
    .option(key: "startingOffsets", value: "earliest") \
    .load()

kafka_df1 = kafka_df.select(from_avro(col('value'), schema_avro).alias("data")) \
    .select(col("data.ad_id"), col("data.timestamp"), col("data.clicks"), col("data.views"), col("data.cost")) \
    .withColumn( colName: 'current_timestamp', current_timestamp())

result_df = kafka_df1.groupBy(window(col("current_timestamp"), windowDuration: "1 minute", slideDuration: "30 seconds"), col("ad_id")) \
    .agg(
    sum("clicks").alias("total_clicks"),
    sum("views").alias("total_views"),
    avg("cost").alias("avg_cost_per_view")
) \
    .select("window.start", "window.end", "ad_id", "total_clicks", "total_views", "avg_cost_per_view")
```

**Fig.: Streaming Code**

```python
def process_batch(batch_df, batch_id):
    max_retries = 3
    retry_count = 0
    connected = False
    while not connected and retry_count < max_retries:
        try:
            # Initialize Cassandra connection
            cluster = Cluster(['cassandra'], port=9042)
            session = cluster.connect('spark')
            connected = True
        except NoHostAvailable:
            retry_count += 1
            print(f"Retry {retry_count}/{max_retries} - Unable to connect to Cassandra. Retrying...")
            if retry_count == max_retries:
                print("Failed to connect to Cassandra after several retries.")
                return

    try:
        # Process each row in the batch
        for row in batch_df.collect():
            ad_id = row['ad_id']
            total_clicks = row['total_clicks']
            total_views = row['total_views']
            avg_cost_per_view = row['avg_cost_per_view']

            # Query Cassandra
            query = SimpleStatement(
                f"SELECT total_clicks, total_views, avg_cost_per_view FROM adsData WHERE ad_id = '{ad_id}'")
            result = session.execute(query).one()

            if result:
                new_total_clicks = result.total_clicks + total_clicks
                new_total_views = result.total_views + total_views
                new_avg_cost_per_view = (
                                            result.avg_cost_per_view * result.total_views + avg_cost_per_view * total_views) / new_total_views

                update_query = SimpleStatement(
                    f"UPDATE adsData SET total_clicks = {new_total_clicks}, total_views = {new_total_views}, avg_cost_per_view = {new_avg_cost_per_view} WHERE ad_id = '{ad_id}'")
                session.execute(update_query)
```

**Fig.: Streaming Code**

```
            session.execute(update_query)
        else:
            insert_query = SimpleStatement(
                f"INSERT INTO adsData (ad_id, total_clicks, total_views, avg_cost_per_view) VALUES ('{ad_id}', {total_clicks}, {total_views}, {avg_cost_per_view})")
            session.execute(insert_query)

        # Clean up
        session.shutdown()
        cluster.shutdown()
    except Exception as e:
        print(f"Error processing batch: {e}")


query = result_df.writeStream \
    .foreachBatch(process_batch) \
    .outputMode("update") \
    .start()

query.awaitTermination()
```

**Fig.: Streaming Code**

## Step 4:

After that using CQL I created keyspace and then create table in which I will be writing data make sure to match the schema of the table with the datatypes. after running the job I checked table for the entry of data and also checked the topic which I was producing messages simultaneously checked the log of the container for running spark application.



**Fig.: Message producing in bytes format on confluent with running spark application container.**

**Fig.: Result Table in Cassandra**

## Challenges:

1. setup the environment using docker.
2. It was not able to connect with spark-avro and spark-sql-kafka (Lots of time spent on trying to get the correct jar/jdbc drivers).
3. Lot of time spent on finding logic to process the data in cassandra for each row.