

- Using a shorter time window resulted in relatively quicker but more random results. With a shorter time window, the results were more varied.

```

25     }
26     producer.send(TOPIC, event)
27     print(f"Sent event: {event}")
28     time.sleep(random.uniform(0.5, 2.0)) # random interval
29

```

PROBLEMS 9 OUTPUT DEBUG CONSOLE TERMINAL PORTS AZURE

```

Sent event: {'song_id': 101, 'timestamp': 1742242742.736314, 'region': 'US', 'action': 'skip'}
Sent event: {'song_id': 303, 'timestamp': 1742242743.6386163, 'region': 'US', 'action': 'like'}
Sent event: {'song_id': 404, 'timestamp': 1742242744.233173, 'region': 'APAC', 'action': 'skip'}
Sent event: {'song_id': 404, 'timestamp': 1742242746.04255, 'region': 'APAC', 'action': 'play'}
Sent event: {'song_id': 101, 'timestamp': 1742242747.185981, 'region': 'EU', 'action': 'skip'}
Sent event: {'song_id': 505, 'timestamp': 1742242748.0372932, 'region': 'EU', 'action': 'like'}
Sent event: {'song_id': 101, 'timestamp': 1742242749.2616951, 'region': 'US', 'action': 'skip'}
Sent event: {'song_id': 404, 'timestamp': 1742242750.4120479, 'region': 'APAC', 'action': 'like'}
Sent event: {'song_id': 303, 'timestamp': 1742242751.9350548, 'region': 'US', 'action': 'play'}
Sent event: {'song_id': 303, 'timestamp': 1742242752.8594189, 'region': 'US', 'action': 'play'}
Sent event: {'song_id': 202, 'timestamp': 1742242754.449979, 'region': 'US', 'action': 'play'}

```

- Updated 'now_trending'.py:

```

o # now_trending.py
o
o from pyspark.sql import SparkSession
o from pyspark.sql.functions import col, count, desc, window,
  from_json, sum as spark_sum
o from pyspark.sql.functions import desc
o from pyspark.sql.types import StructType, StructField,
  StringType, DoubleType
o from pyspark.sql.types import TimestampType
o from pyspark.sql.window import Window
o from pyspark.sql.functions import row_number
o
o # 1) Create SparkSession
o spark = SparkSession.builder \
o     .appName("NowTrendingSongs") \
o     .getOrCreate()
o

```

```

o spark.sparkContext.setLogLevel("WARN")
o
o # 2) Read from Kafka
o kafka_df = spark.readStream \
o     .format("kafka") \
o     .option("kafka.bootstrap.servers", "localhost:9092") \
o     .option("subscribe", "music_events") \
o     .option("startingOffsets", "latest") \
o     .load()
o
o # 3) Parse the JSON 'value' from Kafka
o schema = StructType([
o     StructField("song_id", StringType(), True),
o     StructField("timestamp", DoubleType(), True), # or we can
o     store as Double
o     StructField("region", StringType(), True),
o     StructField("action", StringType(), True)
o ])
o
o json_df = kafka_df.selectExpr("CAST(value AS STRING) as
o     json_str")
o
o parsed_df = json_df.select(from_json(col("json_str"),
o     schema).alias("data"))
o events_df = parsed_df.select("data.*")
o
o # Convert timestamp double -> actual timestamp if we want
o     event time
o # But for simplicity, let's do a processing-time approach
o # If you want event-time windows, do:
o # events_df = events_df.withColumn("event_time",
o     (col("timestamp") * 1000).cast(TimestampType()))
o
o # 4) Filter only "play" events
o plays_df = events_df.filter(col("action") == "play")
o

```

```

o # 5) Group by region + 5-minute processing time window
o # We'll do a simple processing-time window using
  current_timestamp
o # Alternatively, you can do event-time with a column if you
  convert 'timestamp' to a Spark timestamp
o from pyspark.sql.functions import current_timestamp
o
o windowed_df = plays_df \
o     .groupBy(
o         window(current_timestamp(), "1 minute"), #
processing-time window
o         col("region"),
o         col("song_id")
o     ) \
o     .count()
o
o # 6) Use foreachBatch to do rank-based top N logic each
  micro-batch
o def process_batch(batch_df, batch_id):
o     """
o     This function is called for each micro-batch. We treat
    'batch_df' as a normal batch DataFrame.
o     We'll rank by 'count' within each region & window and pick
    top 3 (or 5, or 100).
o     """
o     if batch_df.rdd.isEmpty():
o         print("No data in this batch.")
o         return
o
o     # We'll partition by region + the 'window' column
o     w = Window.partitionBy("region",
"window").orderBy(desc("count"))
o
o     ranked_df = batch_df.withColumn("rn",
row_number().over(w)) \
o         .filter(col("rn") <= 3)

```

```

o
o     # Show the top songs for each region + 5-min window
o     print(f"=== Batch: {batch_id} ===")
o     ranked_df.show(truncate=False)
o
o # 7) Write Stream with foreachBatch
o query = windowed_df \
o     .writeStream \
o     .outputMode("update") \
o     .foreachBatch(process_batch) \
o     .trigger(processingTime='5 seconds') \
o     .start()
o
o skip_ratio_df = (
o     events_df.groupBy("song_id")
o     .agg(
o         count("*").alias("total_events"),
o         count(col("action")).alias("total_skips"),
o     )
o     .withColumn("skip_ratio", col("total_skips") /
o         col("total_events"))
o )
o
o query = (
o     skip_ratio_df.writeStream
o     .outputMode("complete")
o     .format("console")
o     .start()
o )
o
o query.awaitTermination()
o

```

- Updated music_producer.py:

```

o # producer.py

```

```

o import time
o import random
o import json
o from kafka import KafkaProducer
o
o # pip install kafka-python
o
o TOPIC = "music_events"
o producer = KafkaProducer(
o     bootstrap_servers="localhost:9092",
o     value_serializer=lambda v: json.dumps(v).encode('utf-8')
o )
o
o songs = [101, 202, 303, 404, 505]    # sample song IDs
o regions = ["US", "EU", "APAC"]
o actions = ["play", "skip", "like"]
o
o while True:
o     event = {
o         "song_id": random.choice(songs),
o         "timestamp": time.time(),
o         "region": random.choice(regions),
o         "action": random.choice(actions)    # or skip, etc.
o     }
o     producer.send(TOPIC, event)
o     print(f"Sent event: {event}")
o     time.sleep(random.uniform(0.5, 2.0))    # random interval

```