# Assignment 4

**Shao-Yu Huang 018195554**

**Q1.** **Design a Kafka streaming data processing pipeline using the NYC Taxi Trip dataset (https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page). Review the NYC taxi trip data to understand its structures and contents. Answer the following questions about the Kafka streaming data processing pipeline.**

1. **Pipeline Design: Plan a Kafka pipeline to stream and process the data. You MUST include a data processing step in this stream data processing pipeline using Kafka, meaning you cannot directly output the raw data from the consumer. Write a short summary explaining your design this Kafka streaming data processing pipeline in details and how this pipeline could be used in real-world applications.**
   Dataset: 2025(January) Yellow Taxi Trip Records
   Title: Calculates the average fare every 30 minutes
   Data Pipeline detail:
   - Producer: Reads from PARQUET file and sends each taxi trip record to the Kafka Topic "taxi_trip_data," it also will do some data cleaning, such as drop rows that the pickup time is later than the drop-off time and the fare amount that is below or equal to 0. Each record contains fare_amount, pickup_datetime, dropoff_datetime, and passenger_count.
   - Kafka Topic "taxi_trip_data": This raw data topic stores all the incoming taxi trip data.
   - Consumer: Consume data from Topic, filter it for necessary attributes, and compute the rolling average fare for each 30-minute window. The result (average fare) is then printed to the console.
   - Output: The average fare per 30 minutes is either printed in real-time or sent to another topic for monitoring or storing results in a database for future analysis.

   This pipeline enables real-time monitoring of taxi fares, helping businesses track pricing trends, monitor demand, and optimize the fleet. It could also be used to analyze fare anomalies or sudden spikes in pricing, allowing for quick responses or optimizations.

2. **Implementation: Use Python and Kafka to implement your pipeline. Submit ALL the Python scripts you developed.**
   - Producer.py:

```python
1.  from kafka import KafkaProducer
2.  import pandas as pd
3.  import pyarrow.parquet as pq
4.  import json
5.  import time
6.
7.  # Read and clean the data
8.  yellow_trips = pq.read_table('yellow_tripdata_2025-01.parquet')
9.  yellow_trips = yellow_trips.to_pandas()
```

```python
10.
11.  # Data cleaning process
12.  yellow_trips = yellow_trips[yellow_trips['tpep_pickup_datetime'] <=
     yellow_trips['tpep_dropoff_datetime']]
13.  yellow_trips = yellow_trips[yellow_trips['fare_amount'] > 0]
14.
15.  # Kafka Producer Setup
16.  producer = KafkaProducer(
17.  bootstrap_servers=['localhost:9092'], # Kafka broker address
18.  value_serializer=lambda v: json.dumps(v).encode('utf-8') # Serialize records as JSON
19.  )
20.
21.  # Produce data to Kafka topic
22.  for index, record in yellow_trips.iterrows():
23.  message = {
24.  'fare_amount': record['fare_amount'],
25.  'pickup_datetime': record['tpep_pickup_datetime'].strftime('%Y-%m-%d %H:%M:%S'),
26.  'dropoff_datetime': record['tpep_dropoff_datetime'].strftime('%Y-%m-%d %H:%M:%S'),
27.  'passenger_count': record['passenger_count']
28.  }
29.
30.  # Send message to Kafka topic "taxi_trip_data"
31.  producer.send('taxi_trip_data', value=message)
32.
33.  # Print to verify
34.  print(f"Sent: {message}")
35.
36.  # Add a delay between sending records (optional, to simulate real-time processing)
37.  time.sleep(0.1)
38.
39.  producer.flush() # Ensure all messages are sent before closing
40.  producer.close()
```

- Consumer.py:

```python
1.   from kafka import KafkaConsumer
2.   import json
3.   from collections import deque
4.   from datetime import datetime
5.
6.   # Kafka Consumer Setup
7.   consumer = KafkaConsumer(
8.   'taxi_trip_data', # Kafka topic name
9.   bootstrap_servers=['localhost:9092'],
10.  value_deserializer=lambda m: json.loads(m.decode('utf-8')),
11.  group_id='taxi_group'
12.  )
13.
14.  # Initialize a deque to store data for calculating average fare within a 30-minute window
15.  fare_window = deque()
16.  window_duration = 30 * 60 # 30 minutes in seconds
17.
18.  def calculate_average_fare(window_data):
19.  if not window_data:
20.  return 0
```

```
21.  total_fare = sum([record[1] for record in window_data]) # record[1] is fare_amount
22.  return total_fare / len(window_data)
23.
24.  # Consumer behavior: Processing the streamed data from Kafka
25.  for message in consumer:
26.  record = message.value
27.  fare_amount = record['fare_amount']
28.  pickup_datetime = datetime.strptime(record['pickup_datetime'], '%Y-%m-%d %H:%M:%S')
29.
30.  # Add new record to the fare_window
31.  fare_window.append((pickup_datetime, fare_amount))
32.
33.  # Remove records that are older than the 30-minute window
34.  while fare_window and (pickup_datetime - fare_window[0][0]).total_seconds() >
     window_duration:
35.  fare_window.popleft()
36.
37.  # Calculate and display the average fare for the current window
38.  avg_fare = calculate_average_fare(fare_window)
39.  print(f"Current Pickup Time: {pickup_datetime} | Average fare for the last 30 minutes:
     ${avg_fare:.2f}")
```

3. **Results: Use Zoom or Panopto to record a video of your screen demonstrating your implemented Kafka streaming data processing pipeline, like the class demo. Share a public link to your video. No need to upload the video file. Make sure the link works and others can view it.**
   Link: https://youtu.be/KowdLCOrdOg

**Q2.** **For this task, find an open-source dataset in CSV format and design a data processing task using the MapReduce framework. Implement this MapReduce job in PySpark in local environment.**

1. **Submit the ALL python scripts you have developed. Share the link of your dataset.**
   Dataset Link (I only use title.crew.tsv.gz and name.basics.tsv.gz datasets):
   https://developer.imdb.com/non-commercial-datasets/

```
1.   import findspark
2.   findspark.init()
3.   from pyspark.sql import SparkSession
4.   from pyspark.sql.functions import explode, split, col
5.   import os
6.   import requests
7.   import gzip
8.   import shutil
9.   spark = SparkSession.builder \
10.  .config("spark.driver.host", "localhost") \
11.  .config("spark.driver.memory", "4g") \
12.  .config("spark.executor.memory", "4g") \
13.  .appName("mr_prac") \
14.  .getOrCreate()
15.  spark
16.
17.  def download_and_unzip(url, output_dir="."):
```

```python
18.    filename = url.split("/")[-1]
19.    local_gz_path = os.path.join(output_dir, filename)
20.    local_tsv_path = local_gz_path.replace(".gz", "")
21.
22.    # Skip download if already exists
23.    if not os.path.exists(local_gz_path):
24.        print(f"Downloading {filename}...")
25.        response = requests.get(url, stream=True)
26.        with open(local_gz_path, 'wb') as f:
27.            shutil.copyfileobj(response.raw, f)
28.        print(f"Downloaded: {local_gz_path}")
29.    else:
30.        print(f"File already exists: {local_gz_path}")
31.
32.    # Unzip .gz to .tsv
33.    if not os.path.exists(local_tsv_path):
34.        print(f"Unzipping {filename}...")
35.        with gzip.open(local_gz_path, 'rb') as f_in:
36.            with open(local_tsv_path, 'wb') as f_out:
37.                shutil.copyfileobj(f_in, f_out)
38.        print(f"Unzipped to: {local_tsv_path}")
39.    else:
40.        print(f"TSV already exists: {local_tsv_path}")
41.
42.    return local_tsv_path
43.
44. # IMDb dataset URLs
45. urls = [
46.    "https://datasets.imdbws.com/title.crew.tsv.gz",
47.    "https://datasets.imdbws.com/name.basics.tsv.gz"
48. ]
49.
50. # Download and unzip all
51. for url in urls:
52.    download_and_unzip(url)
53.
54.
55. # Load IMDb crew data
56. crew_df = spark.read.csv("title.crew.tsv", sep="\t", header=True, inferSchema=True,
       nullValue="\\N")
57.
58. # Load name.basics for mapping nconst -> primaryName
59. names_df = spark.read.csv("name.basics.tsv", sep="\t", header=True, inferSchema=True,
       nullValue="\\N")
60.
61. # Clean and extract directors (some titles have multiple directors)
62. crew_directors = crew_df.select("tconst", explode(split(col("directors"),
       ",")).alias("director_id")).na.drop()
63.
64. # Count number of titles per director
65. director_counts = crew_directors.groupBy("director_id").count().alias("title_count")
66.
67. # Join with names to get actual director names
68. director_with_names = director_counts.join(
69.    names_df.select("nconst", "primaryName"),
```

```
70.    director_counts["director_id"] == names_df["nconst"],
71.    "inner"
72.    ).select("primaryName", "count").orderBy(col("count").desc())
73.
74.    # Show top 20 directors
75.    director_with_names.show(20, truncate=False)
76.
77.    sc = spark.sparkContext
78.    print(sc.statusTracker().getJobIdsForGroup(None))
79.
80.    # stopping the Spark session
81.    spark.stop()
```
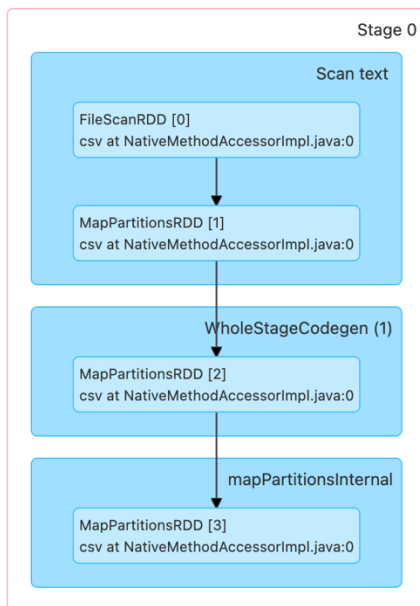
## 2. Provide a screenshot showing the DAG visualization for this MapReduce job.



**Details for Stage 0 (Attempt 0)**

Resource Profile Id: 0
Total Time Across All Tasks: 95 ms
Locality Level Summary: Process local: 1
Input Size / Records: 64.0 KiB / 1
Associated Job Ids: 0

▼ DAG Visualization

**Details for Stage 1 (Attempt 0)**

Resource Profile Id: 0
Total Time Across All Tasks: 19 s
Locality Level Summary: Process local: 8
Input Size / Records: 364.4 MiB / 11584669
Associated Job Ids: 1

▼ DAG Visualization

## Details for Stage 2 (Attempt 0)

**Resource Profile Id:** 0
**Total Time Across All Tasks:** 7 ms
**Locality Level Summary:** Process local: 1
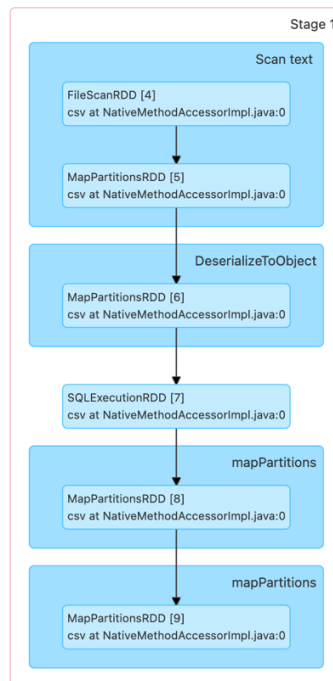**Input Size / Records:** 64.0 KiB / 1
**Associated Job Ids:** 2

▼ DAG Visualization

Stage 2

**Scan text**

FileScanRDD [10]
csv at NativeMethodAccessorImpl.java:0

↓

MapPartitionsRDD [11]
csv at NativeMethodAccessorImpl.java:0

↓

**WholeStageCodegen (1)**

MapPartitionsRDD [12]
csv at NativeMethodAccessorImpl.java:0

↓

**mapPartitionsInternal**

MapPartitionsRDD [13]
csv at NativeMethodAccessorImpl.java:0

## Details for Stage 3 (Attempt 0)

**Resource Profile Id:** 0
**Total Time Across All Tasks:** 24 s
**Locality Level Summary:** Process local: 8
**Input Size / Records:** 840.4 MiB / 14325454
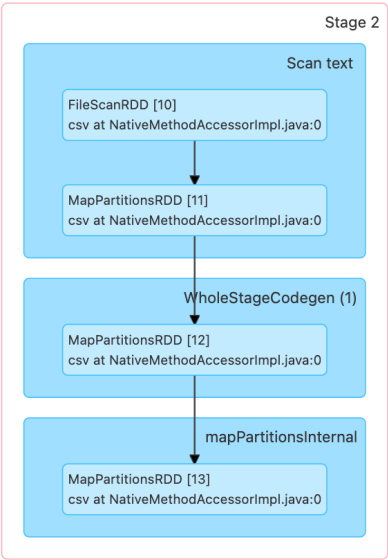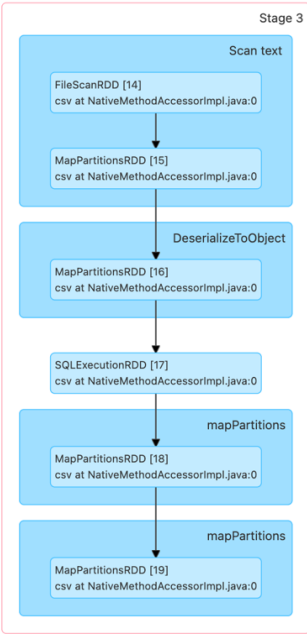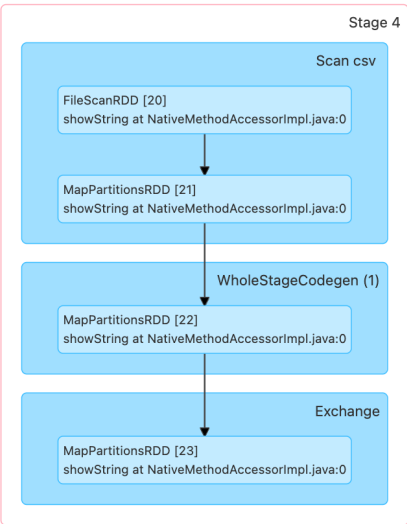**Associated Job Ids:** 3

▼ DAG Visualization

Stage 3

**Scan text**

FileScanRDD [14]
csv at NativeMethodAccessorImpl.java:0

↓

MapPartitionsRDD [15]
csv at NativeMethodAccessorImpl.java:0

↓

**DeserializeToObject**

MapPartitionsRDD [16]
csv at NativeMethodAccessorImpl.java:0

↓

SQLExecutionRDD [17]
csv at NativeMethodAccessorImpl.java:0

↓

**mapPartitions**

MapPartitionsRDD [18]
csv at NativeMethodAccessorImpl.java:0

**mapPartitions**

MapPartitionsRDD [19]
csv at NativeMethodAccessorImpl.java:0

## Details for Stage 4 (Attempt 0)

**Resource Profile Id:** 0
**Total Time Across All Tasks:** 34 s
**Locality Level Summary:** Process local: 8
**Input Size / Records:** 364.4 MiB / 11584668
**Shuffle Write Size / Records:** 18.5 MiB / 1542375
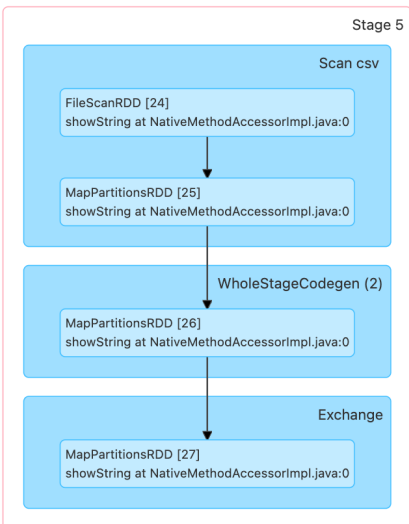**Associated Job Ids:** 4

▼ DAG Visualization

Stage 4

**Scan csv**

FileScanRDD [20]
showString at NativeMethodAccessorImpl.java:0

↓

MapPartitionsRDD [21]
showString at NativeMethodAccessorImpl.java:0

↓

**WholeStageCodegen (1)**

MapPartitionsRDD [22]
showString at NativeMethodAccessorImpl.java:0

↓

**Exchange**

MapPartitionsRDD [23]
showString at NativeMethodAccessorImpl.java:0

## Details for Stage 5 (Attempt 0)

**Resource Profile Id:** 0
**Total Time Across All Tasks:** 1.2 min
**Locality Level Summary:** Process local: 8
**Input Size / Records:** 840.4 MiB / 14325453
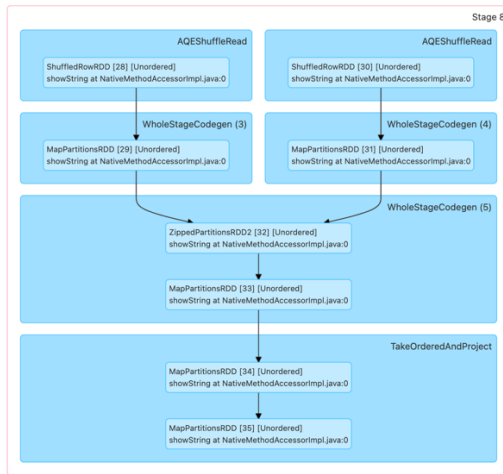**Shuffle Write Size / Records:** 348.4 MiB / 14325453
**Associated Job Ids:** 5

▼ DAG Visualization

Stage 5

**Scan csv**

FileScanRDD [24]
showString at NativeMethodAccessorImpl.java:0

↓

MapPartitionsRDD [25]
showString at NativeMethodAccessorImpl.java:0

↓

**WholeStageCodegen (2)**

MapPartitionsRDD [26]
showString at NativeMethodAccessorImpl.java:0

↓

**Exchange**

MapPartitionsRDD [27]
showString at NativeMethodAccessorImpl.java:0

**Details for Stage 8 (Attempt 0)**

Resource Profile Id: 0
Total Time Across All Tasks: 39 s
Locality Level Summary: Any: 9
Shuffle Read Size / Records: 366.9 MiB / 15867828
Associated Job Ids: 6

▼ DAG Visualization

3. **Provide a summary describing the dataset and the data processing task. Explain how the MapReduce framework was applied to solve the problem.**
   Dataset:
   - "title.crew.tsv": Contains director, writer, and tconst the unique identifier of the title.
   - "name.basics.tsv": Contains metadata for people in the database, including nconst(the unique idnetifier of the name or person), birthYear, deathYear, knownForTitles(titles the person is known for), and primary profession(the top 3 professions of the person).

   The primary goal of the data processing task is to count the number of titles each director has worked on.
   - Map Phase: Each row in the title.crew.tsv file is processed to extract the director IDs. Since titles may have multiple directors (comma-separated), the row is split, and for each director, a tuple (director_id, 1) is emitted — representing one contribution to a title.
   - Shuffle & Sort Phase: Spark automatically groups the intermediate key-value pairs by director_id, preparing them for aggregation.
   - Reduce Phase: The grouped data is reduced by summing the counts, resulting in (director_id, total_titles) pairs.

   To make the output interpretable, the result is joined with name.basics.tsv to map each director_id (nconst) to their primaryName. The final output is a list of director names with the number of titles they have worked on, sorted by count in descending order.

   This MapReduce-based PySpark job efficiently processed thousands of records to output a list of the most active directors in the IMDb database. The highly scalable approach can easily be extended to support more complex queries.