

# **ELEN 6889**

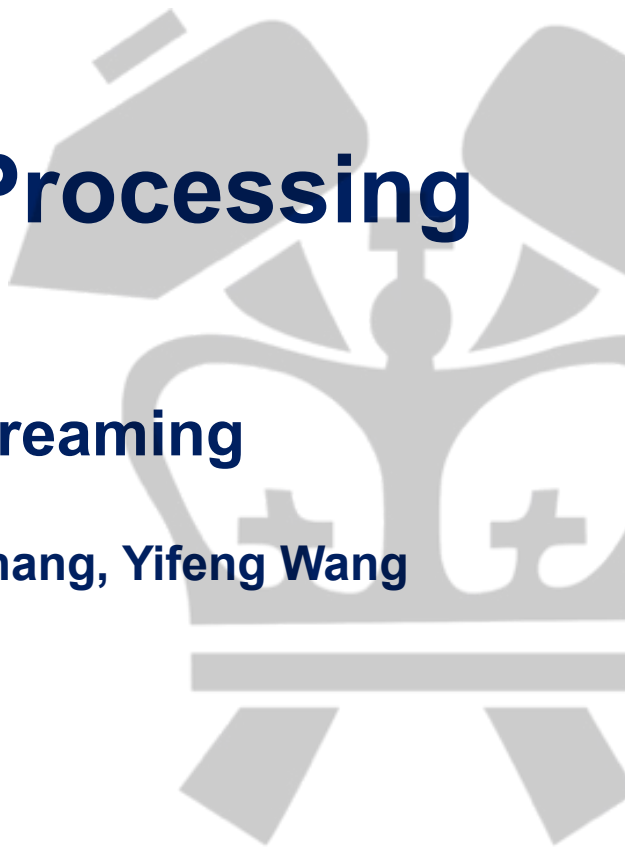
# **Large Scale Data Stream Processing**

## **Final Project Presentation**

## **Topic: Data Analysis in Sport Events Streaming**

**Group members: Weirui Peng, Brian Yao, Jingtian Zhang, Yifeng Wang**

**Date: May 2, 2023**

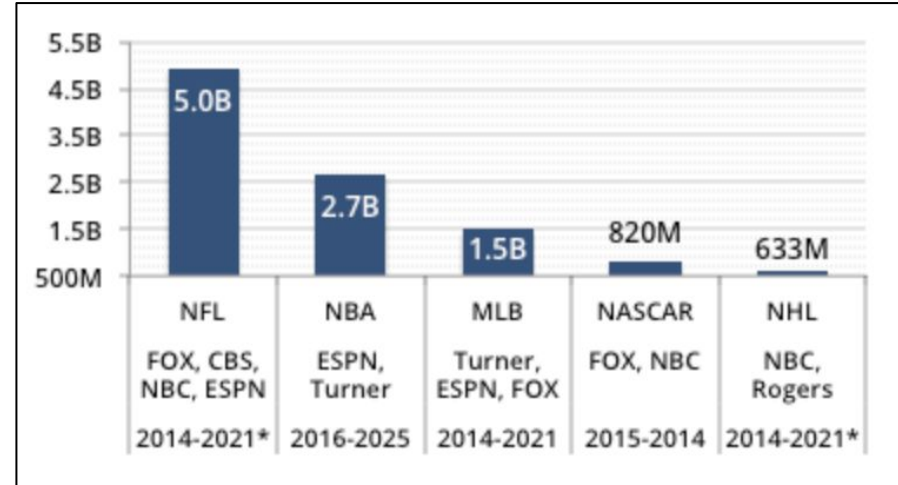


# Contents

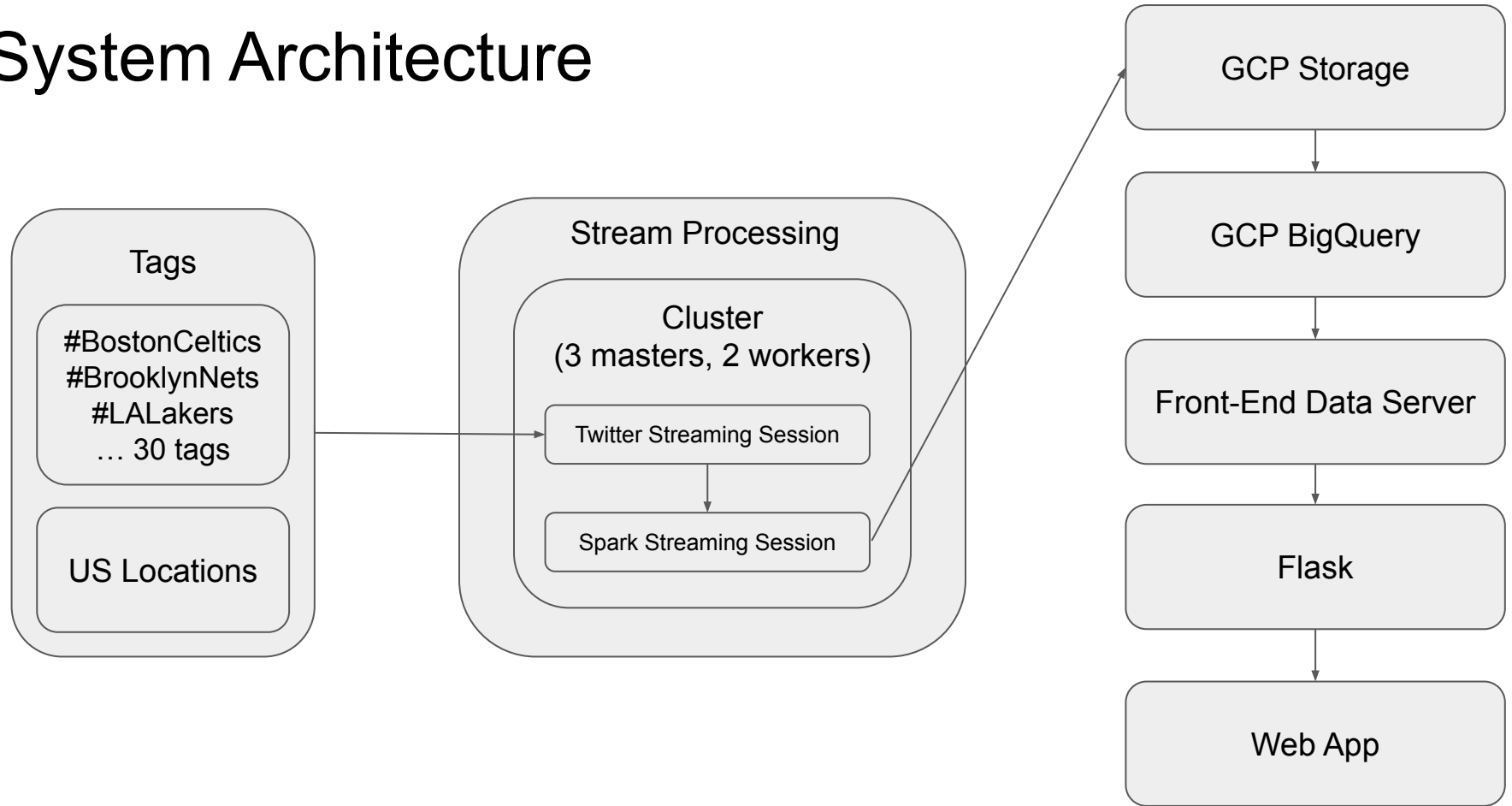
1. Introduction
2. System Architecture
3. Data Selection
4. Analysis Process
5. Optimization
6. Results & Demo

# Introduction

- Sport leagues rely on viewership and ratings for their respective TV deals
- Social media produces data constantly and allows us to look in the lens of a fan to understand how events are perceived
- Real-time streaming allows for real-time strategy adjustment



# System Architecture



# Data Collection

- Set up necessary API keys and access tokens to connect to the Twitter API through Tweepy
- Define search parameters: hashtags, date & time windows, and quantity of tweets to be pulled per hashtag
- Prepare NLTK (Natural Language Toolkit) to be used for sentiment analysis
- For each hashtag, create a CSV file to save the tweet data

# Data Collection

- Parse tweets for query matches. For each tweet:
  - a. Retrieve creation time, text, location (if available), and favorite count
  - b. Analyze sentiment using the NLTK (Natural Language Toolkit)
  - c. Save hashtag, time, text, location, favorite count, and sentiment in a CSV file
- The resulting dataset consists of several CSV files, each containing the tweet data for a single hashtag, including the hashtag, time, text, location, favorite count, and sentiment of each tweet

# Data Collection

- Difficulties
  - a. Twitter API level limit
    - Apply an ELEVATED level twitter developer access
    - Access data in recent 7 days
    - So only collect data during 4/13/2023 - 4/21/2023
  - b. Twitter API tweet limit
    - At most 15,000 tweets within 15 minutes
    - So the program will automatically wait for 15 minutes when facing “429 error”

# Data Collection

```
16 auth = tweepy.OAuth1UserHandler(api_key, api_secret, access_token, access_secret)
17 api = tweepy.API(auth)
18
19 until_date = '2023-04-22'
20 until_time = '23:59:59'
21
22 max_tweets = 200
23 query_date_time = datetime.strptime(until_date + ' ' + until_time, '%Y-%m-%d %H:%M:%S')
24 query_date_time_formatted = query_date_time.strftime('%Y-%m-%d_%H:%M:%S')
25
26 query_hashtags = []
27 # 'ClevelandCavaliers', 'NewYorkKnicks', 'Philadelphia76ers', 'BostonCeltics', 'AtlantaHawks', 'MilwaukeeBucks', 'MiamiHeat', '
28 # '76ers', 'Nets', 'Celtics', 'Bucks', 'Cavaliers', 'Knicks', 'Heat', 'Hawks'
29 # 'Nuggets', 'Grizzlies', 'Kings', 'Suns', 'Clippers', 'Warriors', 'Lakers', 'Timberwolves'
30 # 'PhoenixSuns', 'DenverNuggets', 'MemphisGrizzlies', 'SacramentoKings', 'LAClippers', 'GoldenStateWarriors', 'LosAngelesLaker
```

```
38 for hashtag in query_hashtags:
39     query = '#' + hashtag + ' until:' + query_date_time_formatted
40     path = hashtag + '.csv'
41     with open(path, 'a', encoding='utf-8', newline='') as f:
42         csv_writer = csv.writer(f)
```

```
48
49
50 print('Now print: ' + hashtag)
51 for tweet in tweepy.Cursor(api.search_tweets, q=query, count=max_tweets).items():
52     time = tweet.created_at.strftime('%Y-%m-%d %H:%M:%S')
53     text = tweet.text
54     if tweet.place is not None:
55         location = tweet.place.full_name + ', ' + tweet.place.country
56     else:
57         location = 'No location available.'
58     favorite_count = tweet.favorite_count
59
60 # 分析推文的情感
61 sentiment_scores = sia.polarity_scores(text)
62 sentiment = 0 # 中性情感
63 if sentiment_scores['pos'] > sentiment_scores['neg']:
64     sentiment = 1 # 积极情感
65 elif sentiment_scores['pos'] < sentiment_scores['neg']:
66     sentiment = -1 # 消极情感
67
68 # 将推文数据作为一行写入 CSV 文件
69 csv_writer.writerow([hashtag, time, text, location, favorite_count, sentiment])
70 sum += 1
71
72 print(sum)
73 except Exception as e:
74     print(e)
75     if '429' in str(e):
76         print('Rate limit exceeded. Waiting for 15 minutes...')
77         t.sleep(15 * 60) # 等待 15 分钟
78         continue
```





# Analysis Process

- Read the streaming data from a Google Cloud Storage bucket using `spark.readStream.csv()`, returning a `DataFrame` that represents the data stream
- Filter the streaming data based on non-null values for location, likes, and sentiment, converting -1 values to 0 for column 'likes', as well as filtering out foreign language text; Use Watermark function on the filtered data helps the system handle late data by defining a 1-hour threshold

```
46  filtered_data = (  
47      streaming_data.filter(  
48          "isnotnull(location) AND isnotnull(likes) AND isnotnull(sentiment)")  
49      .withColumn("likes", when(col("likes") == -1, 0).otherwise(col("likes")))  
50      .filter(is_english_udf(col("text")))  
51      .withWatermark("timestamp", "1 hour")  
52  )
```

# Analysis Process

- Perform windowed aggregation on the filtered data using the `groupBy()` and `agg()` methods.
  - The data is grouped by a 30-minute window and slide duration, as well as the location
  - This allows for real-time analysis of the data within each window
  - Select the necessary columns and rename the window start and end columns to provide a more comprehensive output

```
# Aggregate data
aggregated_data = (
    filtered_data.groupBy(
        window("timestamp", "30 minutes", "30 minutes"),
        "location"
    )
    .agg(
        sum("likes").alias("total_likes"),
        sum(when(col("sentiment") == 1, 1).otherwise(0)).alias("count_sentiment_1"),
        sum(when(col("sentiment") == 0, 1).otherwise(0)).alias("count_sentiment_0"),
        sum(when(col("sentiment") == -1, 1).otherwise(0)).alias("count_sentiment_-1"),
        avg("sentiment").alias("average_sentiment")
    )
    .select(
        "location",
        "total_likes",
        "count_sentiment_1",
        "count_sentiment_0",
        "count_sentiment_-1",
        "average_sentiment",
        col("window").getField("start").alias("window_start"),
        col("window").getField("end").alias("window_end")
    )
)
```

# Analysis Process

- The aggregated data is written to a checkpoint folder and new data will be appended to the existing output data
- The final output data is read into a DataFrame and sorted based on the window\_start column in ascending order

```
# Write aggregated data to temporary folder
query = (
    aggregated_data.writeStream
        .outputMode("append")
        .format("parquet")
        .option("path", "gs://dataproc-staging-us-central1-300388429198-ocbfs3ky/outputdata/")
        .option("checkpointLocation", "gs://dataproc-staging-us-central1-300388429198-ocbfs3ky/checkpoint/")
        .start()
)

query.awaitTermination()

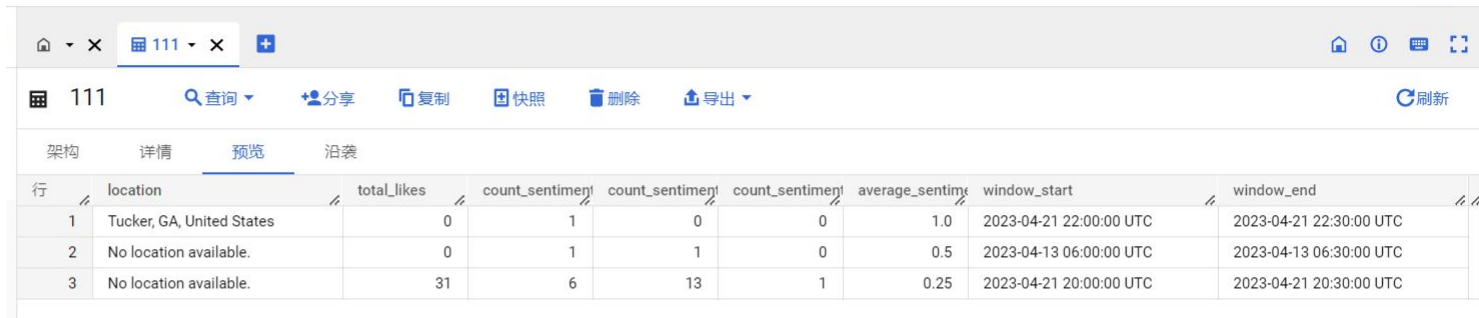
# Read the output from the temporary folder
output_data = spark.read.parquet("gs://dataproc-staging-us-central1-300388429198-ocbfs3ky/outputdata/")

# Sort the output data and write the final result
sorted_output = output_data.sort(asc("window_start"))

sorted_output.write.csv("gs://dataproc-staging-us-central1-300388429198-ocbfs3ky/outputdata/", mode="overwrite")
```

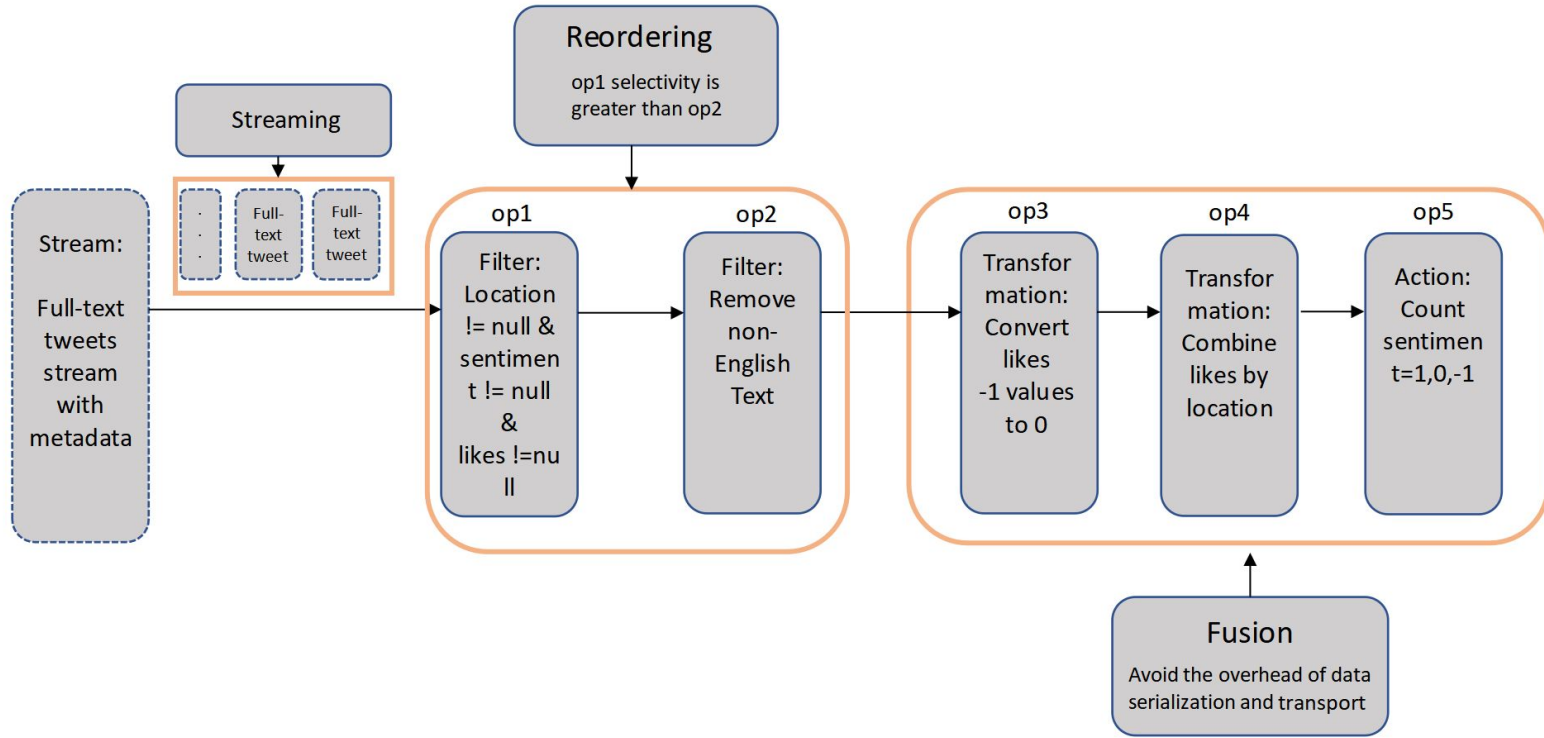
# Analysis Process

- The sorted output data is written as CSV files to a Google Cloud Storage bucket, overwriting any existing data



行	location	total_likes	count_sentimen	count_sentimen	count_sentimen	average_sentime	window_start	window_end
1	Tucker, GA, United States	0	1	0	0	1.0	2023-04-21 22:00:00 UTC	2023-04-21 22:30:00 UTC
2	No location available.	0	1	1	0	0.5	2023-04-13 06:00:00 UTC	2023-04-13 06:30:00 UTC
3	No location available.	31	6	13	1	0.25	2023-04-21 20:00:00 UTC	2023-04-21 20:30:00 UTC

# Optimization



# Results & Demo

- Main results:
  - a. Visualization
  - b. Query powered by BigQuery
- Our website is deployed on GCP:  
<http://34.66.247.178:8111/>

# Results & Demo

- Our plot function is extendable by adding new queries

```
8 def plot_bigquery(start_time, end_time):
9     # create BigQuery client object
10    client = bigquery.Client()
11
12    # define SQL query to retrieve the required data
13    query = f"""
14        SELECT string_field_0 AS location,
15               int64_field_1 AS count_likes,
16               int64_field_2 AS count_sentiment1,
17               int64_field_3 AS count_sentiment0,
18               int64_field_4 AS count_sentiment_1,
19               double_field_5 AS avg_sentiment,
20               timestamp_field_6 AS start_time,
21               timestamp_field_7 AS end_time
22    FROM `NBAdataset.result`
23    WHERE timestamp_field_6 >= TIMESTAMP("{start_time.strftime('%Y-%m-%d %H:%M:%S')} UTC") AND timestamp_field_6 <= TIMESTAMP("{end_time.strftime('%Y-%m-%d %H:%M:%S')}")
24    """
25
26    # execute the query and convert the result to a pandas DataFrame
27    df = client.query(query).to_dataframe()
28
29    path1 = plot_1(df)
30    path2 = plot_2(df)
31    path3 = plot_3(df)
32
33    return path1, path2, path3
```

Thank you for your attention