

Final Project

Professor: Saber Amini

Members: Chen, Kun (8977010)

Kalaiarasan, Abinaya (8953115)

Yalamkayala, Haripriya Naga Neelima. (8939330)

Get NewsAPI token

Registration complete

Your API key is: `c76ec7ae04464220bba2b1af2e170b05`

For help getting started please look at our [getting started guide](#).

We post API status updates and other news on our Twitter feed, so please follow us there if that's important to you:

[Follow @NewsAPlorg](#)

 My account

Start zookeeper & kafka server, then create a new topic: DonaldTrump

```

1 # in zookeeper/
2 bin/zkServer.sh start
3
4 # in confluent/
5 nohup bin/kafka-server-start etc/kafka/server.properties > /dev/null 2>&1 &
6
7 # Create topic: DonaldTrump
8 bin/kafka-topics --create --zookeeper localhost:2181 --partitions 1 --replication-factor 1
--topic DonaldTrump
9 # Show topic describe
10 bin/kafka-topics --describe --zookeeper localhost:2181 --topic DonaldTrump
11

```

Setup a Kafka producer

We use the exactly same source code from midterm project.

1. Get all article from NewsAPI which mention "**Donald Trump**"
2. Pass those article to the consumer via kafka,

Producer

```

1 import requests
2 from confluent_kafka import Producer
3 from dotenv import load_dotenv
4 import os
5 import time
6 import hashlib
7
8 load_dotenv()
9 api_key = os.getenv('NEWS_API_KEY')
10
11 news_url = f'https://newsapi.org/v2/everything?q=Donald&q=Trump&sortBy=popularity&apiKey={api_key}'
12
13 producer = Producer({'bootstrap.servers': 'localhost:9092'})
14 topic = 'DonaldTrump'
15 processed_ids = set()
16
17 def delivery_report(err, msg):
18     if err is not None:
19         print(f'Message sent failed: {err}')
20     else:
21         print(f'Message Send: {msg.topic()} [{msg.partition()}]')
22
23 while True:
24     try:
25         response = requests.get(news_url)
26         response.raise_for_status()
27         news_data = response.json()
28
29         for article in news_data['articles']:

```

```

30         url = article.get('url', 'Null')
31         title = article.get('title', 'No title')
32         publishedAt = article.get('publishedAt', '0')
33
34         source = article['source'].get('name', 'Null')
35         author = article.get('author', 'Null')
36         description = article.get('description', 'No description')
37         urlToImage = article.get('urlToImage', 'Null')
38         content = article.get('content', 'Null')
39
40         article_id = hashlib.md5(f'{title}{publishedAt}{url}'.encode('utf-
8')).hexdigest()
41         article['article_id'] = article_id
42
43         message = str(source) + '|' + str(author) + '|' + str(description) + '|' +
str(urlToImage) + '|' + str(content)
44
45         if article_id not in processed_ids:
46             processed_ids.add(article_id)
47             producer.produce(topic, message.encode('utf-8'), callback=delivery_report)
48
49             producer.flush()
50         else:
51             continue
52
53         print("Wait Next Pull ... ")
54         time.sleep(10)
55
56     except requests.exceptions.RequestException as e:
57         print(f'Request Failed: {e}')
58         time.sleep(60)
59     except Exception as e:
60         print(f'Error: {e}')
61         break

```

Spark Streaming

In consumer side:

1. Get message from Kafka
2. Process the data
3. Aggregate the data
4. Save the data to HDFS

```

1  from pyspark.sql import SparkSession
2  from pyspark.sql.functions import *
3  from pyspark.sql.types import *
4
5  spark = SparkSession.builder.appName("NewsStreamingConsumer") \
6      .config("spark.sql.shuffle.partitions", "2") \
7      .getOrCreate()

```

```

8
9 spark.sparkContext.setLogLevel("WARN")
10
11 raw_df = spark.readStream.format("kafka") \
12     .option("kafka.bootstrap.servers", "localhost:9092") \
13     .option("subscribe", "DonaldTrump") \
14     .option("startingOffsets", "latest") \
15     .load()
16
17 string_df = raw_df.selectExpr("CAST(value AS STRING) as value")
18
19 split_df = string_df.select(
20     split(col("value"), "\|").getItem(0).alias("source"),
21     split(col("value"), "\|").getItem(1).alias("author"),
22     split(col("value"), "\|").getItem(2).alias("description"),
23     split(col("value"), "\|").getItem(3).alias("urlToImage"),
24     split(col("value"), "\|").getItem(4).alias("content"),
25     current_timestamp().alias("timestamp")
26 )
27
28 aggregated_df = split_df \
29     .withWatermark("timestamp", "1 minute") \
30     .groupBy(window(col("timestamp"), "1 minute"), col("source")).count()
31
32 query = aggregated_df.select(col("window").start.alias("window_start"), \
33     col("window").end.alias("window_end"), col("source"), col("count")) \
34     .writeStream.outputMode("append").format("csv") \
35     .option("path", "/BigData/news_data/") \
36     .option("checkpointLocation", "/BigData/news_checkpoint") \
37     .start()
38
39 query.awaitTermination()

```

Start producer

Start Spark Streaming

```

gibdata8451@cluster-e861-m:~/final_test$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.4.1 spark_consumer_hdfs.py
:: loading settings :: url = jar:file:/usr/lib/spark/jars/ivy-2.5.1.jar!/org/apache/ivy/core/settings/ivysettings.xml
Ivy Default Cache set to: /home/gibdata8451/.ivy2/cache
The jars for the packages stored in: /home/gibdata8451/.ivy2/jars
org.apache.spark#spark-sql-kafka-0-10_2.12 added as a dependency
:: resolving dependencies :: org.apache.spark#spark-submit-parent-c3acf098-8db7-4a8a-bdd8-bfd295a805bf;1.0
  confs: [default]
    found org.apache.spark#spark-sql-kafka-0-10_2.12;3.4.1 in central
    found org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.4.1 in central
    found org.apache.kafka#kafka-clients;3.3.2 in central
    found org.lz4#lz4-java;1.8.0 in central
    found org.xerial.snappy#snappy-java;1.1.10.1 in central
    found org.slf4j#slf4j-api;2.0.6 in central
    found org.apache.hadoop#hadoop-client-runtime;3.3.4 in central
    found org.apache.hadoop#hadoop-client-api;3.3.4 in central
    found commons-logging#commons-logging;1.1.3 in central
    found com.google.code.findbugs#jsr305;3.0.0 in central
    found org.apache.commons#commons-pool2;2.11.1 in central
:: resolution report :: resolve 866ms :: artifacts dl 32ms
  :: modules in use:
    com.google.code.findbugs#jsr305;3.0.0 from central in [default]
    commons-logging#commons-logging;1.1.3 from central in [default]
    org.apache.commons#commons-pool2;2.11.1 from central in [default]
    org.apache.hadoop#hadoop-client-api;3.3.4 from central in [default]
    org.apache.hadoop#hadoop-client-runtime;3.3.4 from central in [default]
    org.apache.kafka#kafka-clients;3.3.2 from central in [default]
    org.apache.spark#spark-sql-kafka-0-10_2.12;3.4.1 from central in [default]
    org.apache.spark#spark-token-provider-kafka-0-10_2.12;3.4.1 from central in [default]
    org.lz4#lz4-java;1.8.0 from central in [default]
    org.slf4j#slf4j-api;2.0.6 from central in [default]
    org.xerial.snappy#snappy-java;1.1.10.1 from central in [default]
-----
|           |         modules      ||   artifacts   |
|   conf     | number| search|dwnlded|evicted|| number|dwnlded|
-----|       default | 11  | 0  | 0  | 0  || 11  | 0  |
-----
:: retrieving :: org.apache.spark#spark-submit-parent-c3acf098-8db7-4a8a-bdd8-bfd295a805bf
  confs: [default]
  0 artifacts copied, 11 already retrieved (0kB/18ms)
25/04/19 10:26:58 INFO SparkEnv: Registering MapOutputTracker
25/04/19 10:26:58 INFO SparkEnv: Registering BlockManagerMaster
25/04/19 10:26:58 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
25/04/19 10:26:58 INFO SparkEnv: Registering OutputCommitCoordinator
25/04/19 10:26:58 INFO MetricsConfig: Loaded properties from hadoop-metrics2.properties

```

Data Save at HDFS



```

gibdata8451@cluster-e861-m:~/final_test$ hadoop fs -ls /BigData/news_data/
Found 7 items
drwxr-xr-x  - gibdata8451 hadoop          0 2025-04-19 10:39 /BigData/news_data/_spark_metadata
-rw-r--r--  1 gibdata8451 hadoop      128 2025-04-19 10:39 /BigData/news_data/part-00000-222e5d88-f203-49b5-aab8-b0ab9a164e54-c000.csv
-rw-r--r--  1 gibdata8451 hadoop          0 2025-04-19 10:27 /BigData/news_data/part-00000-53dbc3af-34fa-4c7a-9164-d29ee6dea0e5-c000.csv
-rw-r--r--  1 gibdata8451 hadoop      250 2025-04-19 10:39 /BigData/news_data/part-00000-5ba43b6c-5b00-4ble-94e7-b7646eb155fc-c000.csv
-rw-r--r--  1 gibdata8451 hadoop          0 2025-04-19 10:38 /BigData/news_data/part-00000-f7ad800e-368d-4dd2-84a6-121d029539ba-c000.csv
-rw-r--r--  1 gibdata8451 hadoop      59 2025-04-19 10:39 /BigData/news_data/part-00001-a3862b77-fc93-41c2-a9d6-97e4d4b8009f-c000.csv
-rw-r--r--  1 gibdata8451 hadoop     124 2025-04-19 10:39 /BigData/news_data/part-00001-e9b1cb1d-9e13-4055-a475-eb11bf14fd42-c000.csv

```

Create Hive Table

```

1 CREATE EXTERNAL TABLE news_aggregated (
2   window_start TIMESTAMP,
3   window_end TIMESTAMP,
4   source STRING,
5   count INT
6 )
7 ROW FORMAT DELIMITED
8 FIELDS TERMINATED BY ','
9 STORED AS TEXTFILE
10 LOCATION '/BigData/news_data';

```

Query the Hive Table

```
1 SELECT * FROM news_aggregated LIMIT 10;
```



SSH-in-browser

```
hive> SELECT * FROM news_aggregated LIMIT 10;
OK
NULL      NULL      The Verge        52
NULL      NULL      Gizmodo.com     18
NULL      NULL      Wired          16
NULL      NULL      The Verge        52
NULL      NULL      The Verge        52
NULL      NULL      Gizmodo.com     18
NULL      NULL      Wired          16
NULL      NULL      Wired          16
NULL      NULL      Gizmodo.com     18
Time taken: 1.871 seconds, Fetched: 9 row(s)
```