

COURSE WORK: BIG DATA ANALYTICS TECHNOLOGY

MSc Data Science: Coventry University UK (**2024.2 BATCH**)

Name: P D Lolitha Lakshan Weerasinghe

Index Number: COMScDS242P-003

Table of Contents

| | |
|---|-----------|
| TASK 1: IN-DEGREE DISTRIBUTION ANALYSIS USING APACHE HADOOP AND APACHE SPARK | 4 |
| PART 1: IMPLEMENTATION AND PERFORMANCE COMPARISON..... | 4 |
| <i>Hadoop Cluster Implementation Overview</i> | <i>4</i> |
| <i>Hadoop Job Implementation</i> | <i>4</i> |
| <i>Hadoop Performance Tracking.....</i> | <i>4</i> |
| <i>Spark Cluster Implementation Overview.....</i> | <i>5</i> |
| <i>Spark Job Implementation</i> | <i>6</i> |
| <i>Spark Performance Tracking</i> | <i>6</i> |
| <i>In-degree distribution plots.....</i> | <i>8</i> |
| PERFORMANCE METRICS [INCLUDES SOC-LIVEJOURNAL1 DATASET]..... | 9 |
| <i>Execution Duration</i> | <i>9</i> |
| <i>Memory Usage</i> | <i>10</i> |
| <i>CPU utilization</i> | <i>13</i> |
| <i>Network Utilization.....</i> | <i>15</i> |
| <i>Disk Utilization</i> | <i>20</i> |
| <i>Performance Conclusion</i> | <i>23</i> |
| PART 2: SCALABILITY AND OPTIMIZATION ANALYSIS | 23 |
| <i>Optimizations</i> | <i>23</i> |
| <i>Analysis.....</i> | <i>24</i> |
| TASK 2: REAL-TIME IOT SENSOR DATA ANALYTICS USING APACHE KAFKA | 26 |
| PART 1: SETUP AND ENVIRONMENT CONFIGURATION | 26 |
| PART 2: DATA SOURCE AND PREPROCESSING..... | 29 |
| PART 3: STREAMING DATA PROCESSING AND ANALYSIS..... | 29 |
| 3.1 <i>Data Ingestion</i> | <i>29</i> |
| PART 4: VISUALIZATION AND REPORTING | 33 |
| TASK 3: GRAPH DATABASES USING NEO4J AND DOCKER OBJECTIVE | 37 |
| PART 1 — UNDERSTANDING GRAPH DATABASES | 37 |
| <i>Fraud Detection In Financial Services.....</i> | <i>37</i> |

| | |
|---|-----------|
| <i>Food Discovery</i> | 38 |
| <i>Contextual Search & Recommendations In Ecommerce</i> | 38 |
| PART 2 — IMPLEMENTATION WITH NEO4J (DOCKER SETUP) | 40 |
| TASK 4: WATERMARKS IN REAL-TIME STREAM PROCESSING USING APACHE KAFKA AND APACHE FLINK (DOCKER SETUP) | 47 |
| PART 1 — UNDERSTANDING WATERMARKS | 47 |
| PART 2 — IMPLEMENTATION USING DOCKER (KAFKA + FLINK)..... | 51 |
| <i>Step 1 — Environment Setup</i> | 51 |
| <i>Step 2 — Data and Kafka Topics</i> | 52 |
| <i>Step 3 — Flink Stream Processing</i> | 53 |
| <i>Step 4 — Scaling Experiment</i> | 62 |
| TASK 5 | 69 |

Source Code

- GitHub Repo : <https://github.com/lolitha-lakshan-1/bigdata-assignment>
- NIBM One Drive : https://nibm-my.sharepoint.com/:u:/g/personal/comscds242p-003_student_nibm_lk/IQDsD-L-e362Q6pPy5K0RvwMAQl4TpUnvFGoYv25MMCf-UM?e=whBcDx

Higer Quality Report in Doc Format : https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/COMScDS242P-003_Report.docx

Task 1: In-Degree Distribution Analysis using Apache Hadoop and Apache Spark

Part 1: Implementation and Performance Comparison

Hadoop Cluster Implementation Overview

The Hadoop cluster was implemented using Docker which consists of NameNode, DataNode and YARN services. Resource configuration was done using the configuration files in the `hadoop_config` directory which was mounted into all the containers; therefore all Hadoop daemons and jobs used the same configuration. The `core-site.xml` was used to set the filesystem to `hdfs://namenode:9000` and the temp directory to `/opt/hadoop/tmp` and `hdfs-site.xml` was used to set the HDFS replication factor to `dfs.replication = 1`, which matches a single-DataNode setup. The `yarn-site.xml` was used to point YARN to the limit the ResourceManager and NodeManager resources to 1024 MB and container allocations between 128 MB and 1024 MB and 1 cpu core.

Hadoop Job Implementation

The Hadoop implementation uses two MapReduce jobs to summarize the SNAP edge-list data. The `InDegreeJob1Hadoop` job reads each edge, skips commented lines and counts the number of times each node appears as destination to get that node's in-degree. The `InDegreeJob2Hadoop` job groups those in-degree results to count how many nodes have each degree value, producing a (degree, count) distribution. The `InDegreeMainHadoop` runs both jobs in order and combines the final output into one CSV for reporting.

Source: https://github.com/lolitha-lakshan-1/bigdata-assignment/tree/main/Task_1/hadoop

Hadoop Performance Tracking

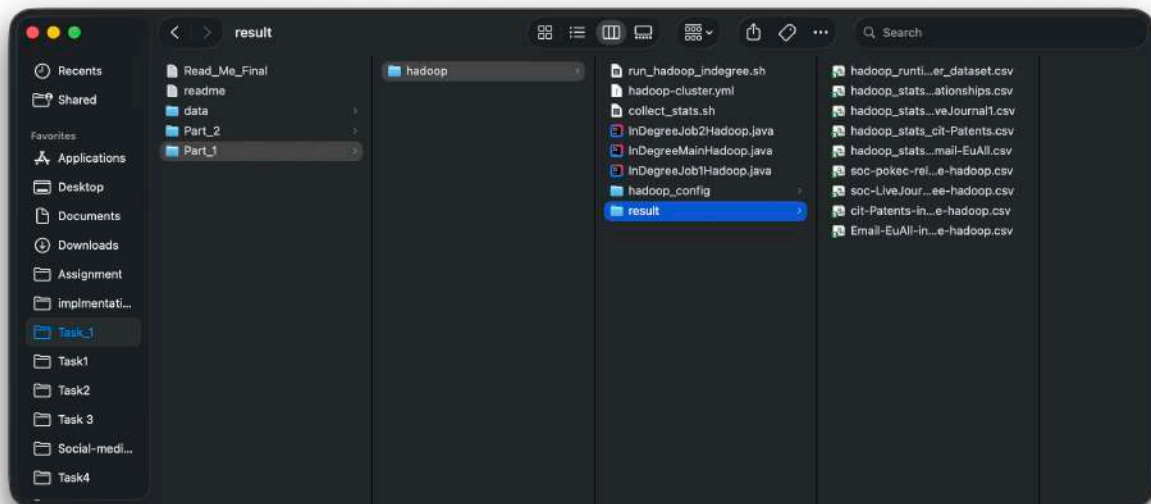
Job performance for each dataset is tracked using a bash script (`run_hadoop_indegree.sh`) which executes the hadoop job per dataset. It logs the job duration and resource usage. For

each dataset, the script first stage the dataset file from DATASET_ROOT into the NameNode and uploaded it to HDFS under /input/<dataset>-hadoop.txt. Then it measures the end-to-end execution time by capturing Unix start/end timestamps around the actual job command and duration is appended to the `hadoop_runtime_per_dataset.csv` file.

In parallel, the CPU, memory usage/percentage, network I/O, disk I/O, and PID metrics are written into a per dataset csv with dataset name `hadoop_stats_<dataset>.csv`. This was done using the executing `docker stats` every 5 seconds for containers. The charts are built using stats from these csv files.

Source : https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_1/hadoop/run_hadoop_indegree.sh

Result : https://github.com/lolitha-lakshan-1/bigdata-assignment/tree/main/Task_1/hadoop/results/PreOptimization



Spark Cluster Implementation Overview

The Spark cluster was implemented using Docker, defined in `spark-cluster.yml`, which consists of a `spark-master` and a `spark-worker` service. The master container starts the Spark Master process, exposes port 7077 for the cluster URL and the worker container starts a Spark Worker process and connects back to the master using port 7077. Both containers

mount to the same host dataset directory “/data”, so the Spark jobs can access the input files across the cluster.

The `run_spark_indegree.sh` script compiles `InDegreeSpark.java`, packages it into `spark-indegree.jar`, copies the JAR into the master container under `/tmp/indegree`, and then runs the job using `spark-submit` on the master with `--master spark://spark-master:7077` and `spark.default.parallelism=4`, passing the dataset name as the input parameter.

Source : https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_1/spark/spark-cluster.yml

Spark Job Implementation

The Spark implementation uses a single Spark job to process the dataset, and it is built using the RDD-based Spark API. Intermediate results are represented as key-value tuples using `JavaPairRDD`. The job reads the input file where each line represents a directed edge as a pair of node IDs. The source node and the destination node. It skips empty lines and comment lines that start with `#`, then extracts the destination node from each edge and emits a tuple `(destinationNode, 1)`. These tuples are aggregated to calculate the in-degree of each node. Thereafter, the job converts those results into `(inDegreeValue, 1)` tuples and aggregates again to count how many nodes share the same in-degree, producing a sorted degree distribution. The final output is written as a single CSV file in the format degree count for reporting.

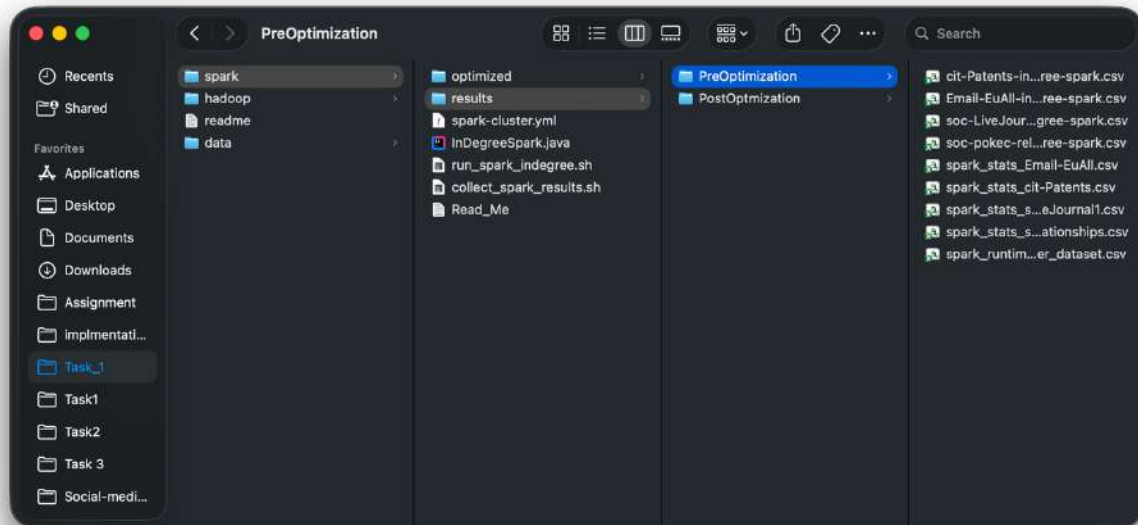
Source: https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_1/spark/InDegreeSpark.java

Spark Performance Tracking

Job performance for each dataset is tracked using the shell script `run_spark_indegree.sh` which executes the Spark job per dataset and logs both duration and resource utilization. For each dataset discovered under `DATASET_ROOT`, the script runs the Spark application using `spark-submit` on the cluster master and measures duration time by capturing Unix

start and end timestamps around the spark-submit command then the duration is appended to spark_runtime_per_dataset.csv.

In parallel, the script starts a background loop that executes every 5 seconds using docker stats command and writes CPU, memory usage/percentage, network I/O, disk I/O, and PID metrics into a per-dataset CSV named spark_stats_<dataset>.csv



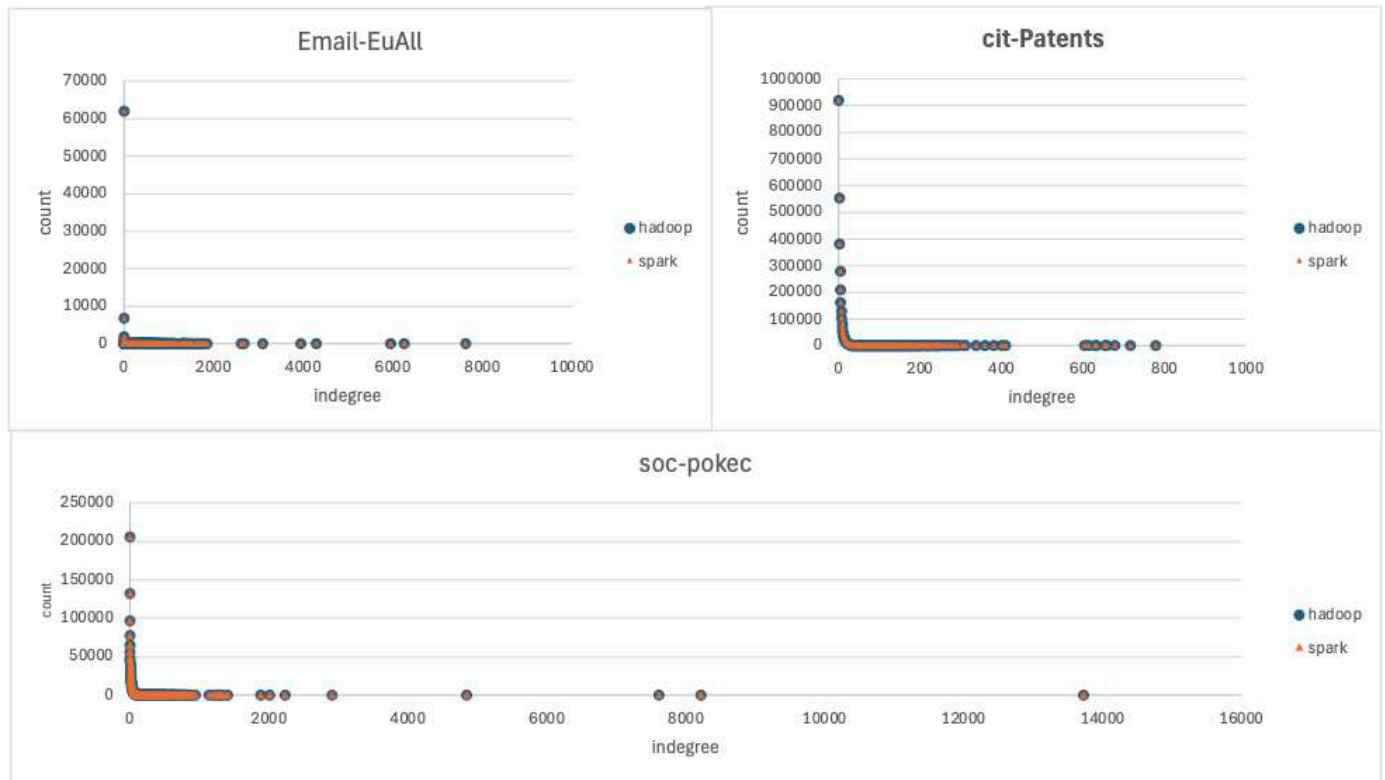
Test Bench

MacBook Air M1 with 8GB RAM and 8-core CPU

In-degree distribution plots

Email-EuAll , cit-Patents and soc-pokec-relationships datasets were executed on the Spark and Hadoop clusters. Their sizes are respectively 5MB, 280MB and 423MB.

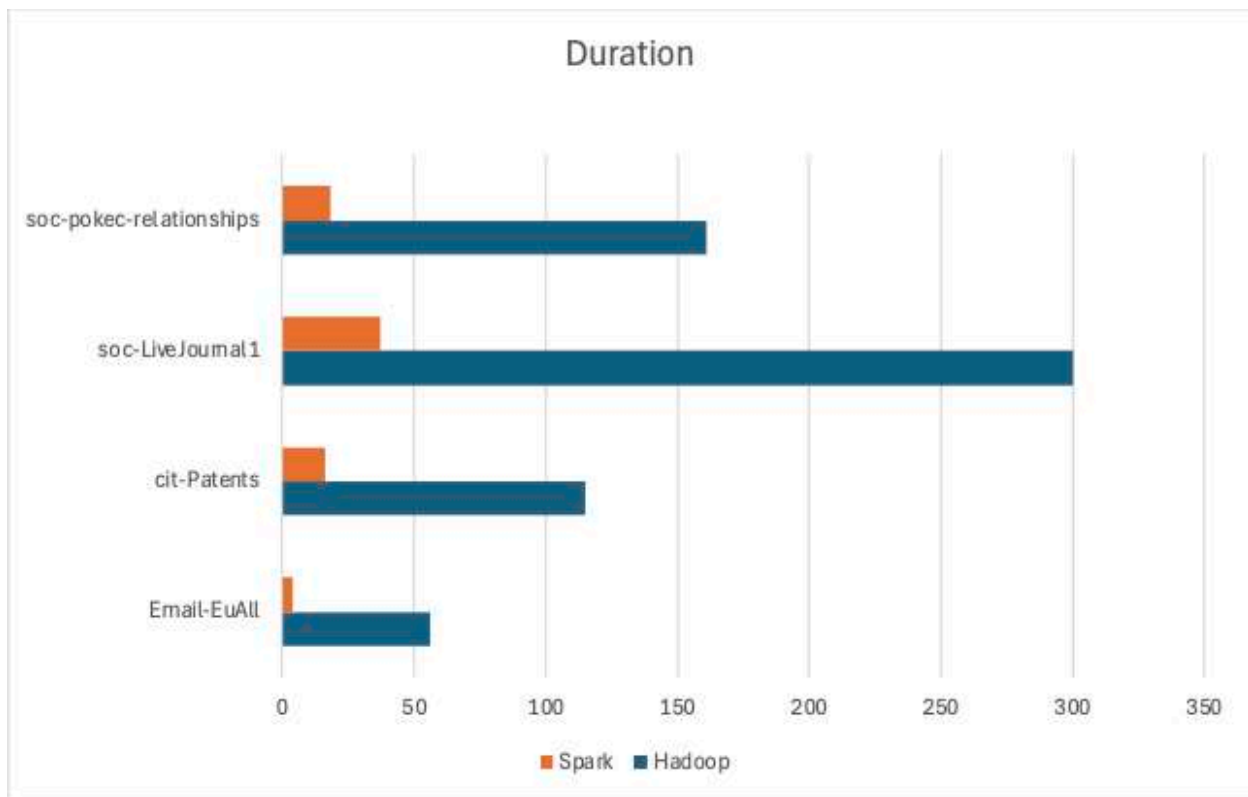
According to the below plots, Hadoop and Spark plots perfectly overlaps for each other for the datasets which confirms that **results are same for both execution methods.**



Performance Metrics [Includes soc-LiveJournal1 dataset]

Excel Sheets With Raw Data : https://github.com/lolitha-lakshan-1/bigdata-assignment/tree/main/Task_1/comparison/pre_optimization

Execution Duration



| Dataset | Size | Hadoop Time(s) | Spark Time(s) | Speedup (Hadoop/Spark) |
|-------------------------|---------|----------------|---------------|------------------------|
| Email-EuAll | 5 MB | 56 | 4 | 14.00x |
| cit-Patents | 280 MB | 115 | 16 | 7.19x |
| soc-pokec-relationships | 423 MB | 161 | 18 | 8.94x |
| soc-LiveJournal1 | 1.08 GB | 300 | 37 | 8.11x |

- Spark is 7.19x to 14x faster across datasets than Hadoop.
- On the smallest dataset (Email-EuAll, 5MB), Hadoop takes 56 seconds. In contrast, it takes only 115 seconds to process a file 56 times larger (280MB). This indicates that Hadoop has a very high startup cost.

- Spark's processing time for the largest dataset (1.08GB) is only 37 seconds, which is still faster than Hadoop's time for the smallest dataset 56 seconds.

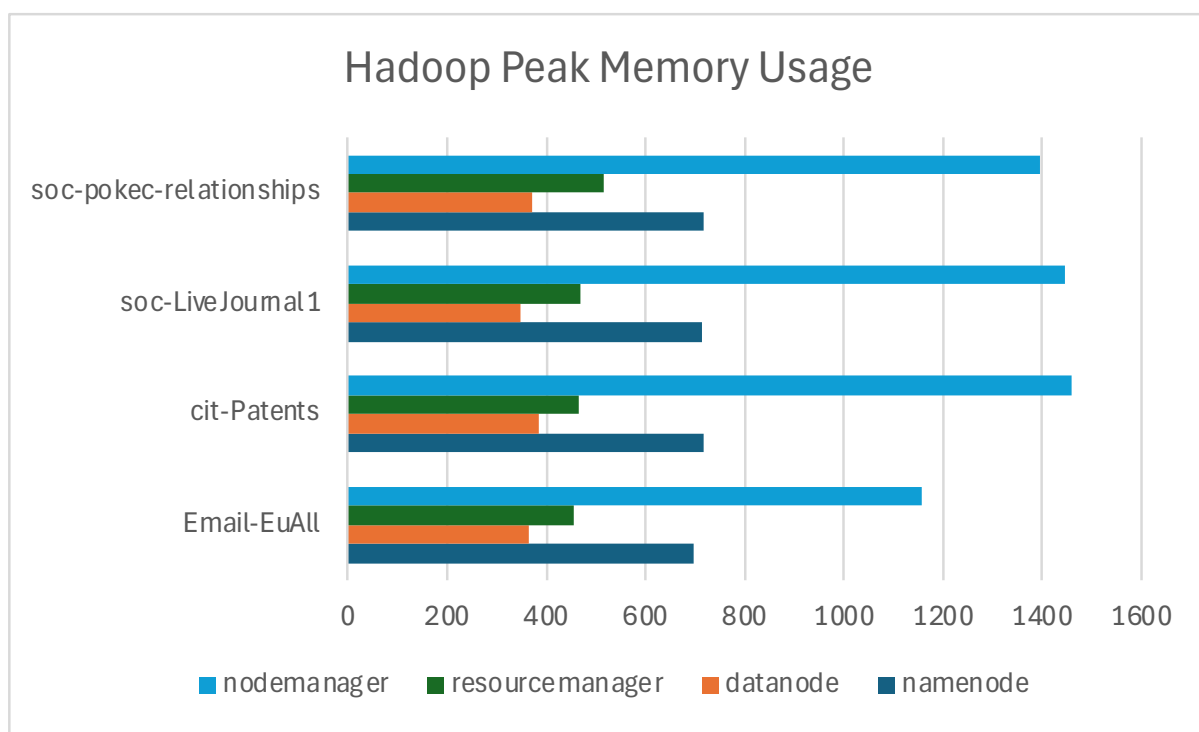
Memory Usage

Peak Memory Usage

It was found out Hadoop components (NameNode, DataNode, ResourceManager) show very little variation in peak memory usage regardless of the dataset size. For example, the NameNode peak only changes by ~20 megabytes between a 5MB and a 1GB file. This suggests Hadoop pre-allocates or maintains a high baseline memory.

Hadoop Peak Memory Usage (MiB)

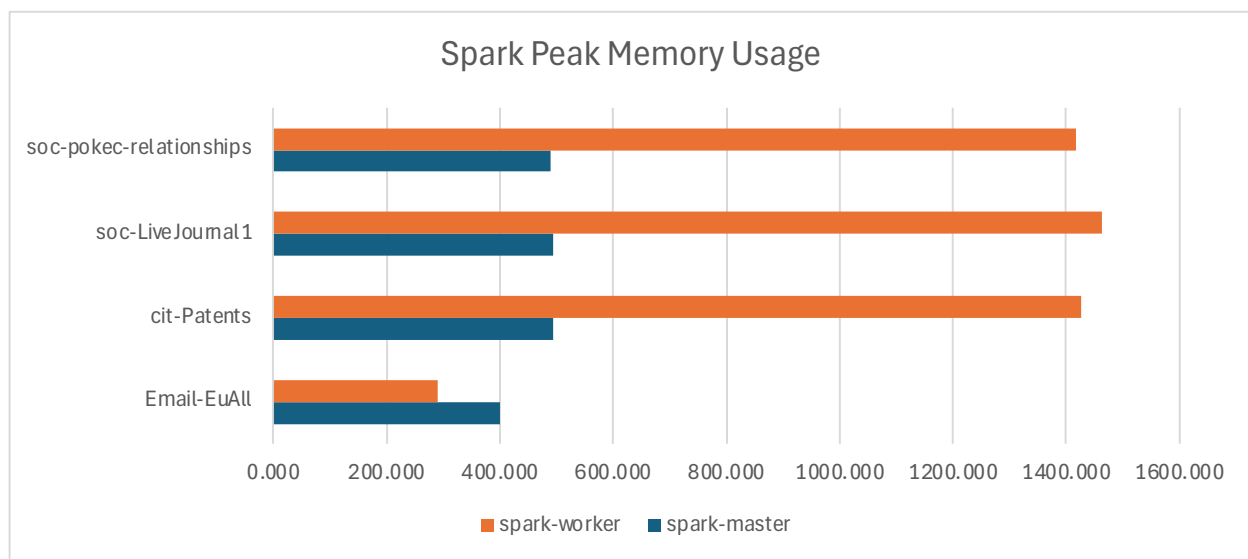
| Dataset | namenode | datanode | resourcemanager | nodemanager |
|-------------------------|----------|----------|-----------------|-------------|
| Email-EuAll | 698.1 | 364.3 | 455.9 | 1158.144 |
| cit-Patents | 718.9 | 383.8 | 464.9 | 1458.176 |
| soc-LiveJournal1 | 715.4 | 348.5 | 470.2 | 1445.888 |
| soc-pokec-relationships | 717.9 | 371.1 | 516.2 | 1395.712 |



- The nodemanager is consistently the largest consumer of memory. It accounts for nearly half of the total peak memory usage for every dataset.
- The namenode and resourcemanager memory peak remains stable regardless of the dataset size.
- The datanode uses the least amount of memory in the cluster, maintaining a steady range between 348.5 and 383.8 megabytes.
- cit-Patents and soc-pokec-relationships show the highest total peak memory usage, largely driven by higher NodeManager requirements. Email-EuAll (the smallest dataset) has the lowest overall memory footprint.

Spark Peak Memory Usage (MiB)

| Dataset | spark-master | spark-worker |
|-------------------------|--------------|--------------|
| Email-EuAll | 399.800 | 291.100 |
| cit-Patents | 494.000 | 1427.456 |
| soc-LiveJournal1 | 494.700 | 1462.272 |
| soc-pokec-relationships | 489.700 | 1416.192 |

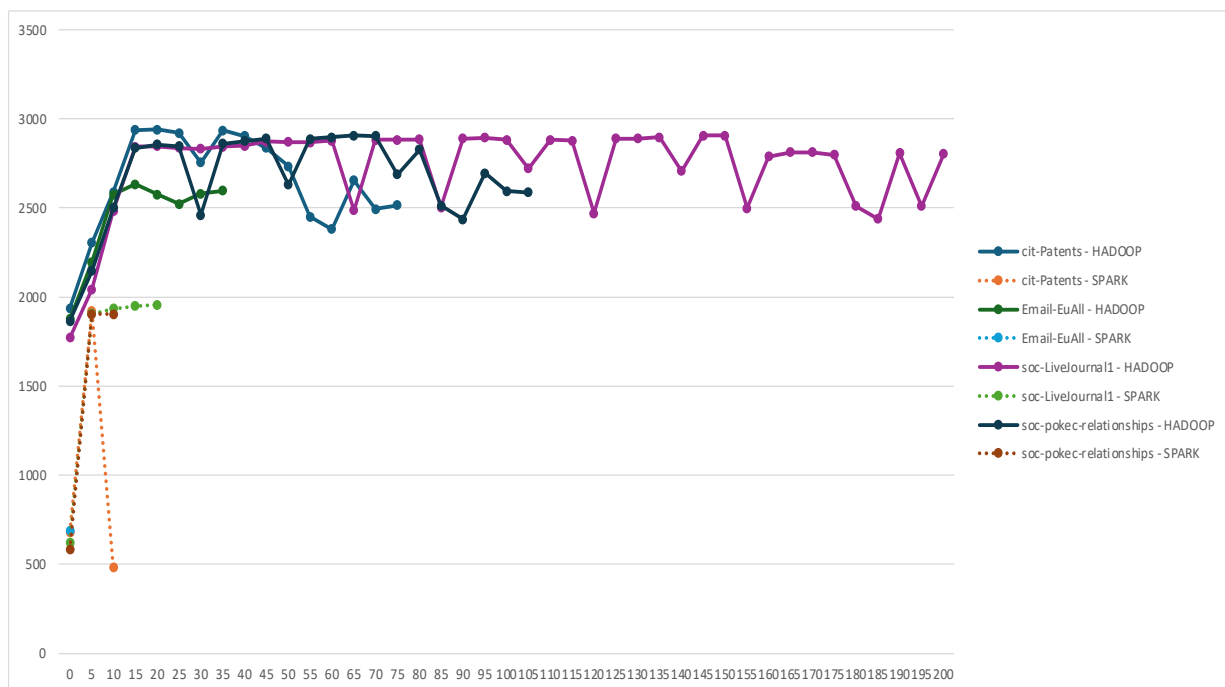


- Spark shows a unique behavior with the smallest dataset, Email-EuAll where spark-master (399.800) uses more memory than the spark-worker (291.100). Suggests

that the administrative overhead of the master node exceeds the actual resources needed for the computation.

- Regardless of whether the dataset is 280MB (cit-Patents) or over 1GB (soc-LiveJournal1), the master stays within a narrow range of 489.7 to 494.7 Mibs.
- The worker memory usage scales sharply as the dataset transitions from small to medium sizes, then plateaus. Memory usage jumps from 291.100 (Email-EuAll) to 1427.456 (cit-Patents).
- However, Spark is much more sensitive to data size. The Spark-Worker peak memory usage jumps from 291.1 for the 5MB file to 1462.2 for the 1GB file.
- Therefore, Spark's memory footprint is small for small tasks and scales up to meet large tasks.

Hadoop vs Spark Memory Usage



For aggregate total of memory usage by all nodes, the chart indicates that Hadoop jobs consumers about 2.5GB - 3GB of memory consistently and the jobs executes for a longer period as well. Meantime Spark jobs peak around 2GB of memory and finish the execution instantly. Therefore, Hadoop maintains a heavy footprint even for the smallest tasks as

depicted above. Furthermore, the solid line for Hadoop depicts a "sawtooth" pattern, where memory fluctuating between 2.5 and 3.0 GB. This indicates shuffle and sort phases of MapReduce. The Spark jobs shoot straight up and straight down. This confirms Spark is loading the graph into memory, processing it instantly.

CPU utilization

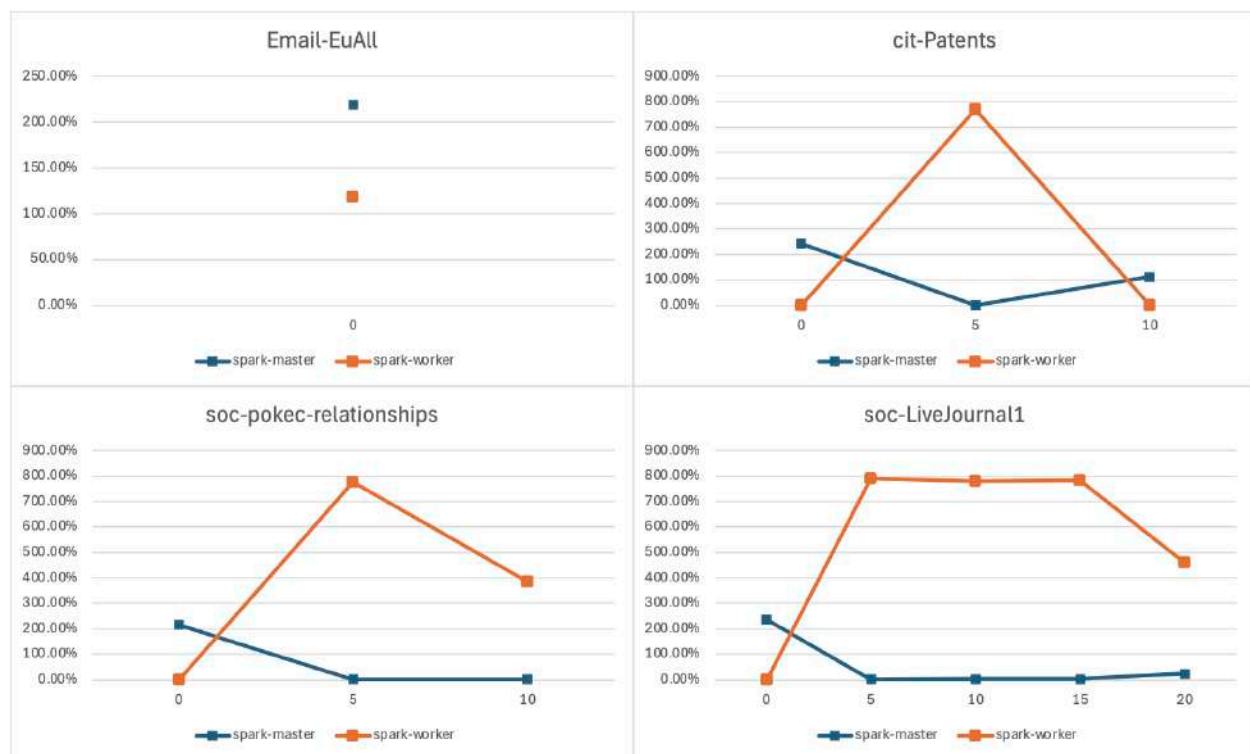
Hadoop CPU Utilization



The values are aggregate across multiple CPU cores (e.g., 200% = 2 full cores utilizing 100% capacity).

- The NodeManager is responsible for launching and managing the actual containers where Map and Reduce tasks run. The high CPU usage confirms that the computation is happening there. Therefore, the green line dominates every chart, staying between 100% and 200%. Also, the "sawtooth" pattern indicates the map reduce steps.
- NameNode (Orange Line) there is a sharp spike at the beginning reaching ~150-200% which immediately drops to near 0%.
- DataNode & ResourceManager (Blue & Light Blue Lines) hover near 0% for the entire duration of the job.
- The chart proves that Hadoop's CPU load is almost entirely on the NodeManagers.
- It seems Hadoop has bottle neck of CPU utilization, only utilizing 2 cores of the available 8 cores. This could be a configuration issue related to docker.

Spark CPU Utilization

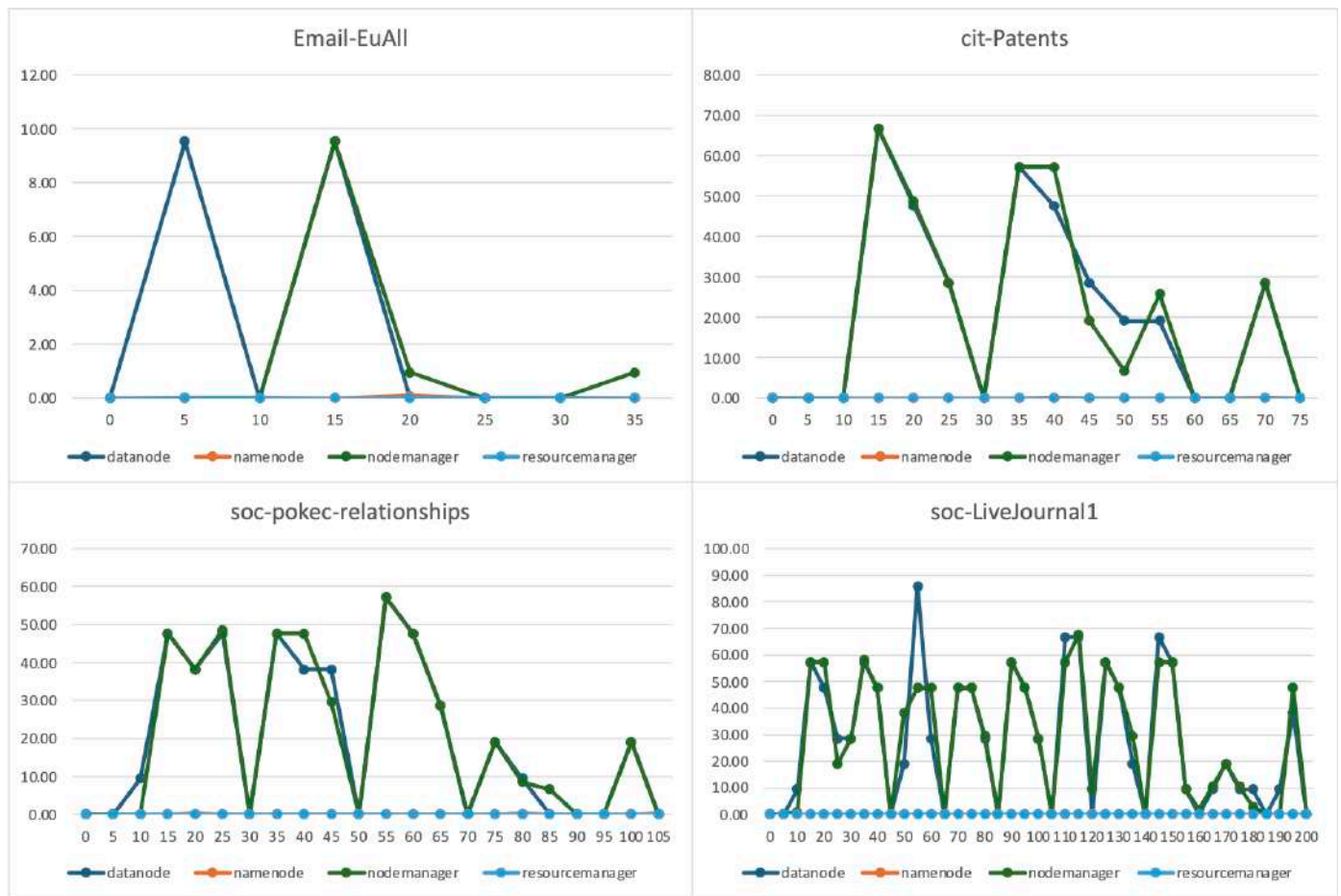


- Unlike Hadoop, which hovered around 100-200% CPU usage, the Spark Worker (Orange) skyrockets to 800% utilization almost immediately for the three larger datasets (cit-Patents, soc-pokec, and soc-LiveJournal1) and 800% peak suggests the worker is fully utilizing the 8-cores in the machine CPU.
- Email-EuAll dataset is so small that the job finishes instantly utilizing only 2 cores.
- cit-Patents & soc-pokec both show a "triangle" shape. The Worker hits 800% at T=5 seconds but is already ramping down by T=10.
- soc-LiveJournal1 this is the only dataset that sustains the load. The Worker hits 800% and stays there for roughly 10-15 seconds. This confirms it is the most computationally intensive job.

Network Utilization

The network utilization values were derived from the `net_io` field, which is recorded as received / sent traffic for each container. First the received side and the sent side were converted into MiB so everything is in one unit. Then a total network value per row was formed as `net_total_mib = net_rx_mib + net_tx_mib`. Because Docker's `net_io` is a cumulative counter, job overhead was measured as the change per 5-second interval for each dataset and each node, the overhead at time t is $\text{net_total_mib_delta}(t) = \text{net_total_mib}(t) - \text{net_total_mib}(t-5)$.

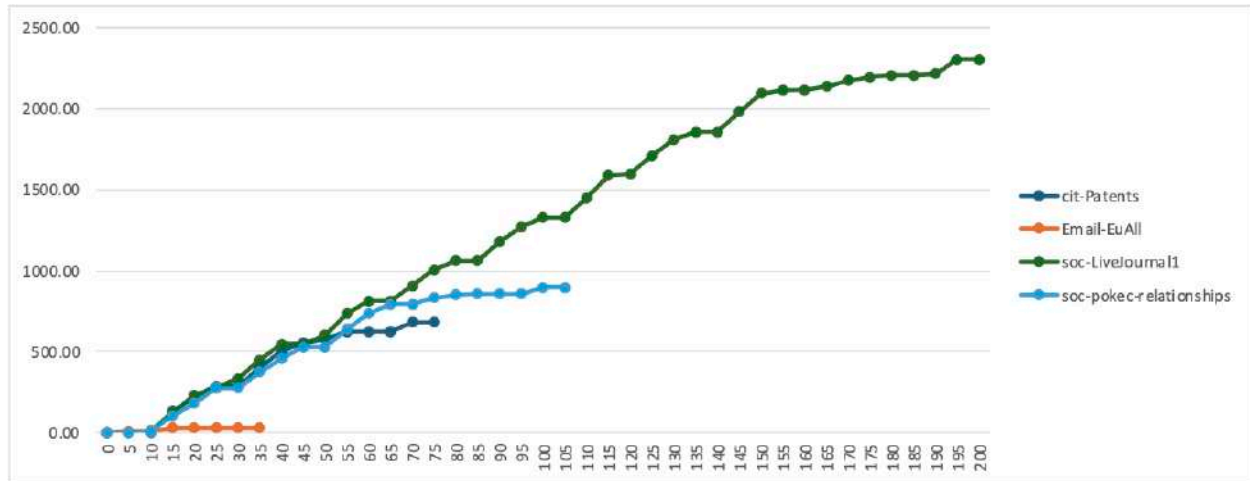
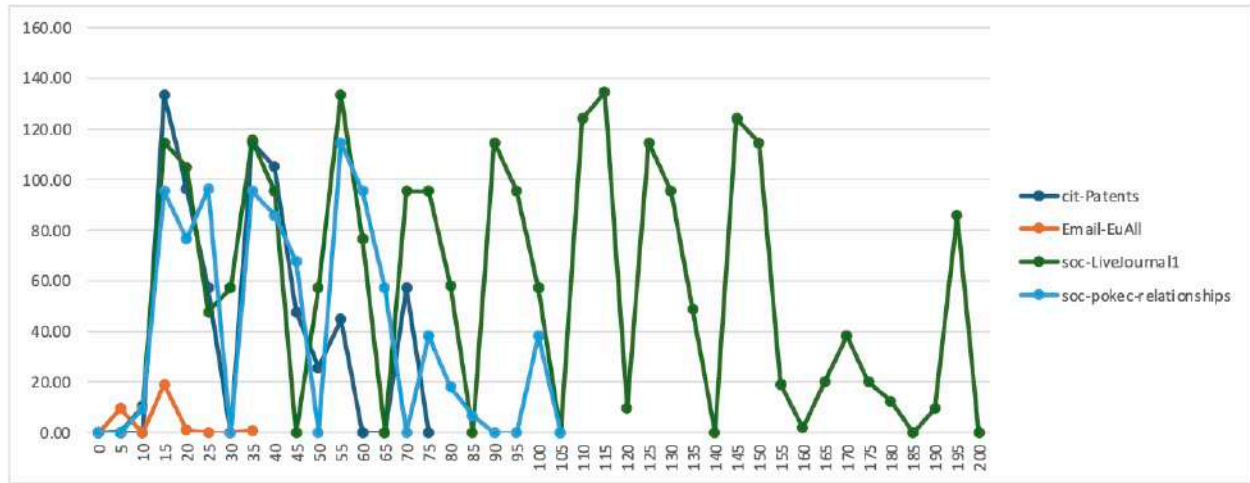
Hadoop Network Utilization



- In all four graphs the DataNode (Blue) and NodeManager (Green) lines move almost perfectly in sync. This confirms that network traffic in Hadoop is driven by data movement between workers. Also, there are sharp peaks followed by drops to near zero which indicates bursty behavior. This would represent the shuffle phase of MapReduce. And the NameNode (Orange) and ResourceManager (Light Blue) remain flat at 0.00 for almost the entire duration across jobs.
- In the smallest job Email-EuAll, the data transfer happens in two quick bursts.
- cit-Patents shows a spike of activity between T=10 and T=60. This indicates a heavy shuffle phase where a large portion of the job execution time was spent just moving data between nodes.

- The soc-pokec dataset shows a series of repeated, medium-intensity spikes. This indicates a constant, rhythmic network shuffling throughout the job.
- soc-LiveJournal1 as the largest dataset, this shows lot of activity. The network is busy constantly from T=15 to T=200, with repeated spikes hitting 60-80 MiB. This suggests the job was network-bound for significant periods, constantly shuffling data back and forth.

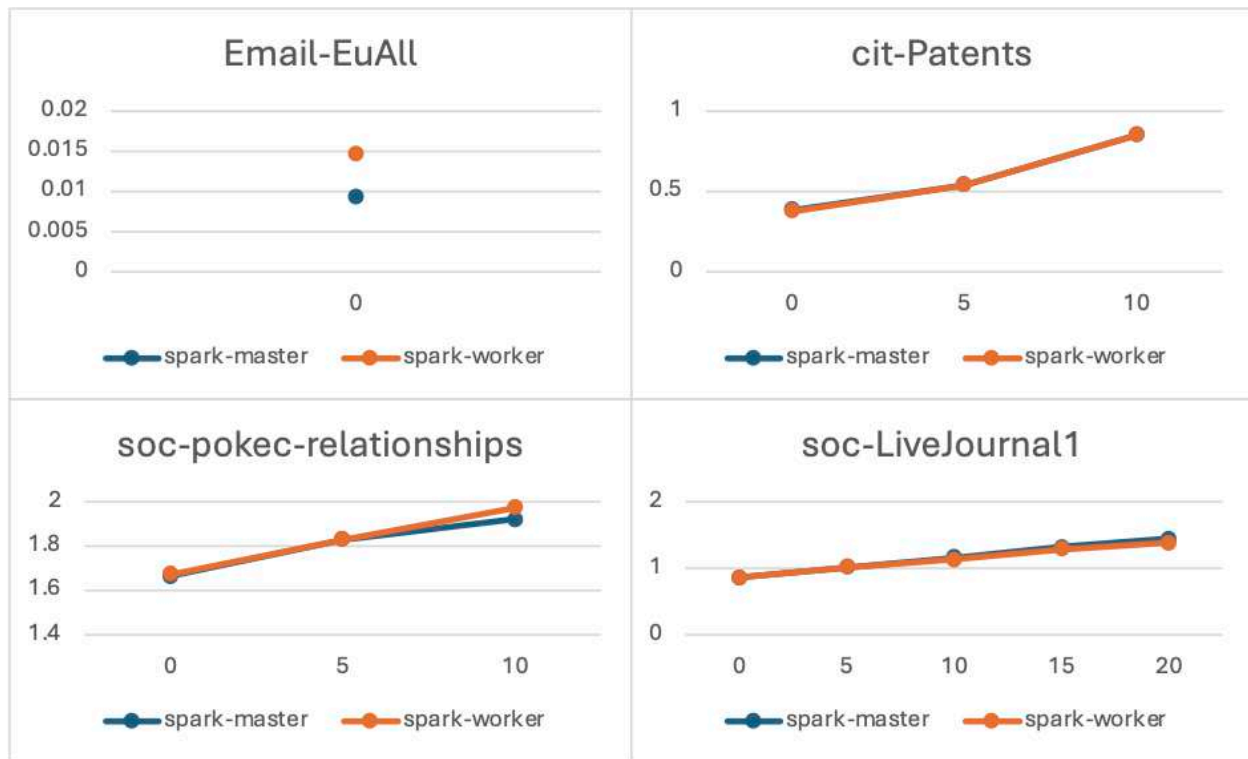
Hadoop Throughput & Cumulative Network IO Usage



The total Network IO profiles demonstrate a direct correlation between data transfer volume and graph edge count. soc-LiveJournal1 (69 million edges) depicts the most extensive network footprint, sustaining a volatile throughput between 10–140 MiB for nearly 200 seconds. The mid-tier datasets such as soc-pokec-relationships (30.6 million edges)

displays a multi-phase pattern with peaks around 115 MiB, while cit-Patents (16.5 million edges) is characterized by a singular, intense initial burst reaching ~135 MiB before dropping off. In contrast, the minimal Email-EuAll (0.42 million edges) requires negligible bandwidth, peaking at only ~20 MiB, confirming that network activity scales aggressively with the volume of edges requiring shuffle.

Spark Network Utilization



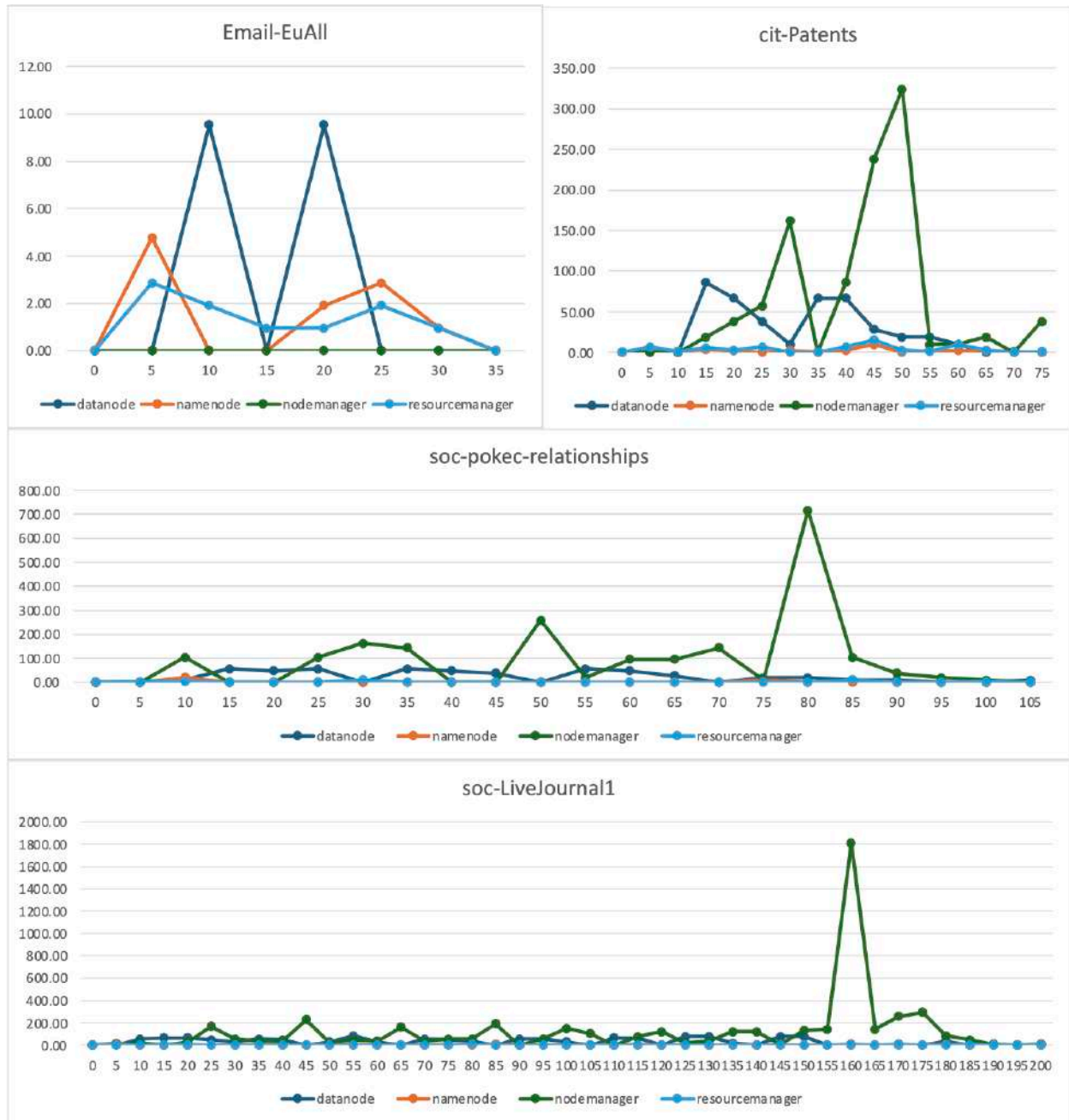
- According to the graphs the Worker nodes shows higher network activity in every dataset than the master nodes. Since Spark Executors resides in Worker nodes, this is as expected as Master nodes act as the cluster manger in a control plane.
- Email-EuAll & cit-Patents shows the lowest network footprint. This Spark job fits mostly in memory, and therefore the data exchange is minimal because the relationships in these datasets are less dense.
- soc-LiveJournal1 & soc-pokec-relationships depicts a jump in Worker network traffic here. Because these are larger social network graphs with high connectivity, Spark must move data between workers.

Hadoop vs Spark Network IO

Based on the findings, there is a massive disparity in network utilization with Hadoop as it generates heavy traffic in the hundreds of MiB, Spark operates almost entirely in the kB range, peaking at only ~1.9 MiB for the same tasks. This ~400x reduction occurs because Hadoop's architecture requires constant data movement between the NodeManager and DataNode for storage and shuffling, whereas Spark processes data in-memory, bypassing the network interface for most of the computation.

Disk Utilization

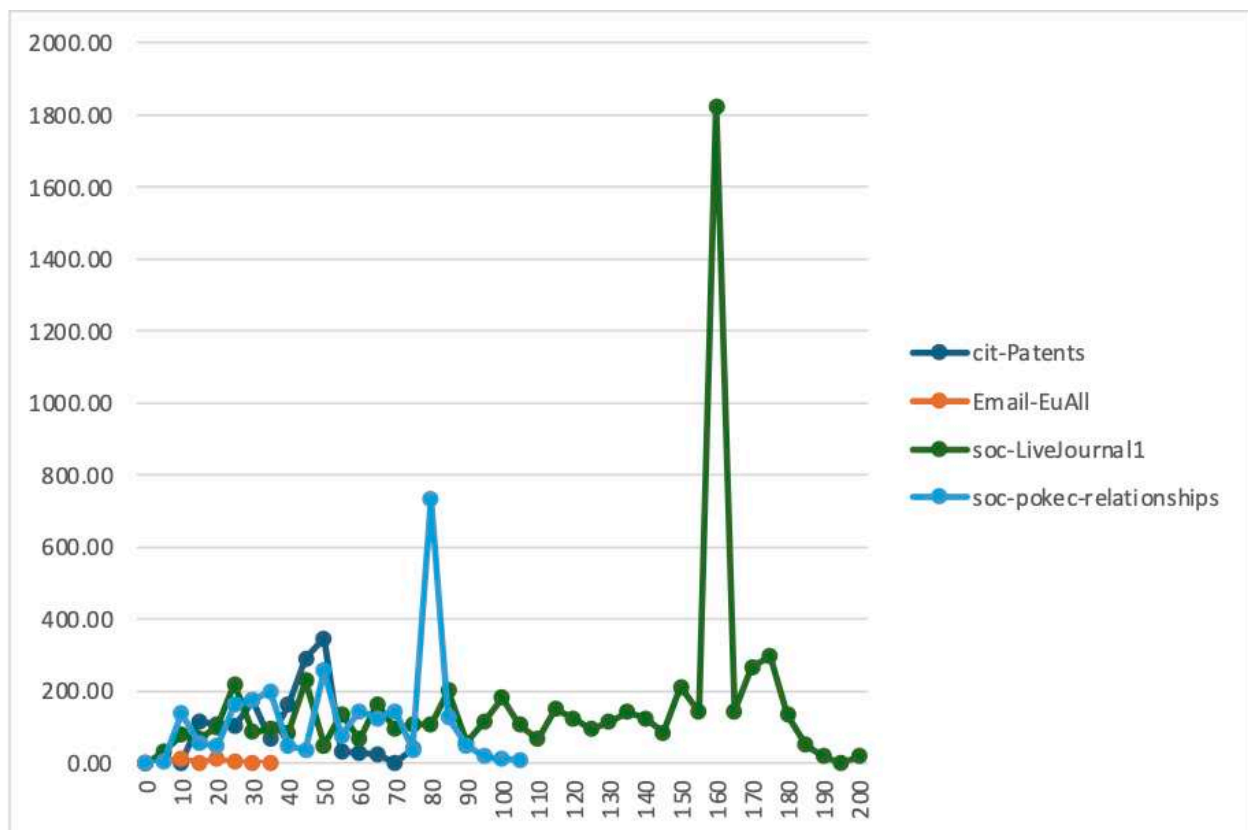
Hadoop Disk Utilization



NameNode and ResourceManager remains idle, using less than 65 MiB of disk per job. It seems NodeManager's activity is the primary driver of disk I/O while DataNode maintains a steady but lower load.

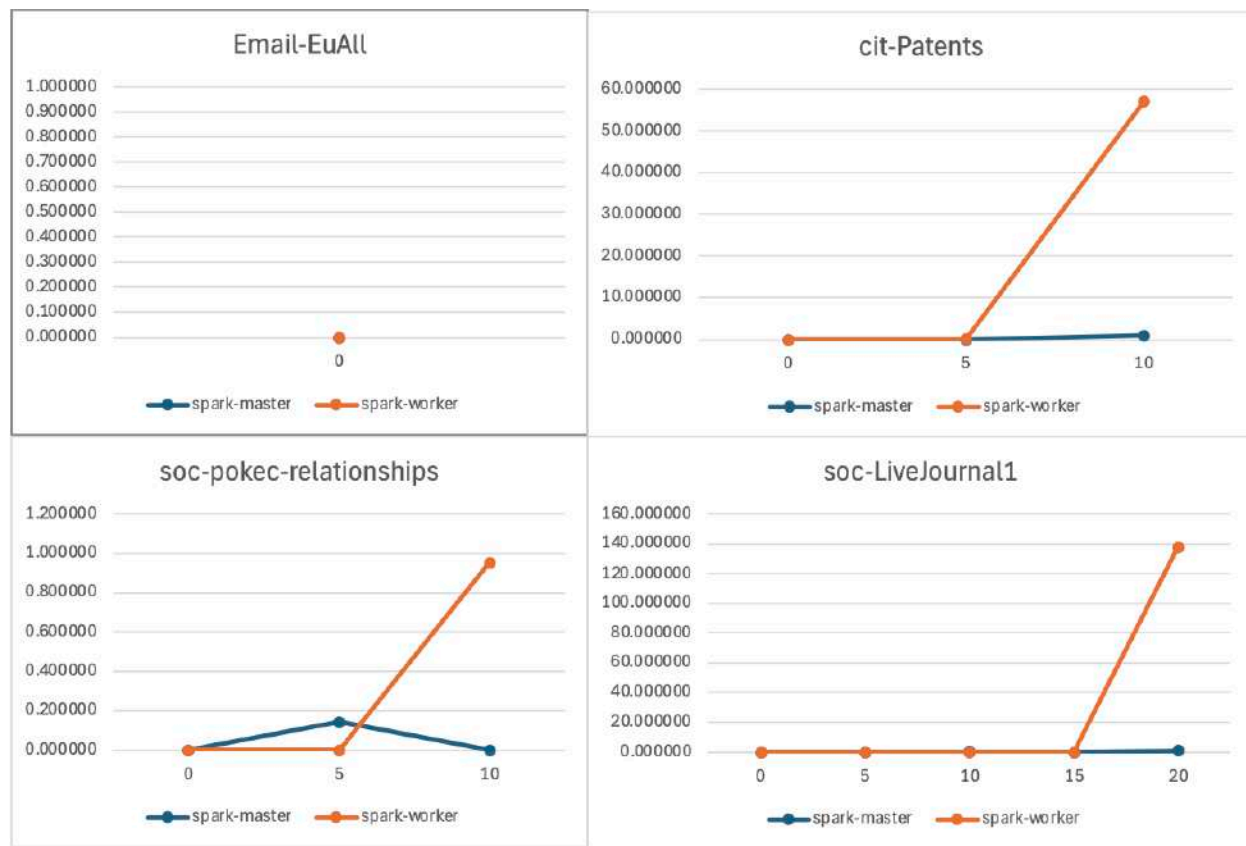
For simple datasets like Email-EuAll, total disk activity is minimal, often staying under 20 MiB. For large dataset like soc-LiveJournal1 causes disk I/O to surge. For the soc-LiveJournal1 dataset, the NodeManager disk io shoots up to a peak of 1,812 MiB in a single interval. This indicates a massive burst of temporary data being written to disk during the shuffle confirming that the "shuffle" creates the highest stress on the disk.

Total Disk IO per Dataset



The chart illustrates the total Disk I/O across four datasets, revealing a correlation between workload complexity and disk IO. The Email-EuAll dataset miniscule footprint. As the workloads scales to cit-Patents and soc-pokec-relationships, disk usage grows linearly, reaching approximately 0.3 GiB and 0.78 GiB. The most significant finding is the disproportionate spike for the soc-LiveJournal1 dataset, which dominates the chart with a peak usage of approximately 1.8GiB which is 2.5x than the soc-pokec-relationships. There seems to be an exponential increase in disk access compared with the dataset sizes.

Spark Disk Utilization



The soc-pokec-relationships dataset traces the highest curve, yet it reaches a peak of only 3.89 MiB at the 10-second mark. The soc-LiveJournal1 workload is even lighter, peaking at 2.82 MiB at the 20-second mark, while cit-Patents registers a maximum of just 1.70 MiB. The Email-EuAll dataset remains a flat baseline near zero.

Spark vs Hadoop Disk IO

There is wide contrast between the disk usage between Hadoop and Spark. While Hadoop's disk usage spiked to 1,822 MiB for the largest workload, Spark's peak activity is over 450 times lower, topping out at less than 4 MiB. This contrast demonstrates how Spark's in-memory architecture fundamentally shifts the execution model, transforming a gigabyte-scale storage workload into a minimal megabyte-scale task.

Performance Conclusion

Hadoops has a disk-centric design and Spark's has a memory-centric architecture. Therefore, Hadoop intermediate data must be written to disk before the next stage begins, resulting in the gigabyte-scale spikes observed in the analysis. Spark, conversely, utilizes an in-memory approach, which allows it to retain intermediate results in RAM. This architectural advantage makes Spark lightweight compared to Hadoop.

Part 2: Scalability and Optimization Analysis

Please refer Part 1 charts and descriptions for performance metrics with soc-LiveJournal1.

Optimizations

Hadoop Optimization

To improve the performance of the Hadoop job, a combiner was introduced. The combiner performs partial aggregation on each mapper prior to the shuffle phase. This reduces the number of intermediate records sent over the network. This decreases the shuffle overhead bandwidth usage and lowers overall job runtime. Since the combiner reduces shuffle volume, so reducers receive and merge map outputs faster and can begin the reduce phase sooner.

| Dataset | Size | Hadoop Time | Optimized Hadoop Time |
|-------------------------|---------|-------------|-----------------------|
| Email-EuAll | 5 MB | 56 | 56 |
| cit-Patents | 280 MB | 115 | 106 |
| soc-pokec-relationships | 423 MB | 161 | 144 |
| soc-LiveJournal1 | 1.08 GB | 300 | 267 |

The optimization had no measurable impact on the smallest dataset. Email-EuAll remained at 56 seconds before and after optimization. For the larger datasets, runtime decreased consistently. cit-Patents improved from 115 seconds to 106 seconds, soc-pokec-relationships improved from 161 seconds to 144 seconds, and soc-LiveJournal1 improved

from 300 seconds to 267 seconds. Overall, the benefit increased with dataset size, showing the optimization is more effective when shuffle volume is higher.

Spark Optimization

The Spark workloads were optimized by increasing the number of partitions Spark uses for parallel operations. This was done using the `spark.default.parallelism=8` property. Initially this was 4 parallel partitions. This will allow Spark to use more CPU cores at the same time making counting and aggregation work is split into more tasks.

| Dataset | Size | Spark Time | Optimized Spark Time |
|-------------------------|-------------|-------------------|-----------------------------|
| Email-EuAll | 5 MB | 4 | 4 |
| cit-Patents | 280 MB | 16 | 16 |
| soc-pokec-relationships | 423 MB | 18 | 18 |
| soc-LiveJournal1 | 1.08 GB | 37 | 36 |

The Spark parallelism change had almost no effect on runtime across the datasets. Email-EuAll stayed at 4 seconds, cit-Patents stayed at 16 seconds, and soc-pokec-relationships stayed at 18 seconds after optimization. Only soc-LiveJournal1 showed a small improvement from 37 seconds to 36 seconds. This indicates the workload was already efficiently parallelized.

Analysis

Hadoop and Spark show different performance patterns because they process data in different ways. Hadoop writes intermediate results to disk between stages and spends more time in shuffle and file I O. Spark keeps more data in memory and reduces repeated disk writes. This gives Spark a large advantage in end-to-end runtime. The Hadoop combiner reduced shuffle volume and improved runtime more on larger datasets, which matches the observed gains.

Spark is better suited for processing large scale graph data in this experiment. It completed all datasets much faster than Hadoop. The Hadoop optimization improved results but Hadoop remained slower, especially on the largest dataset. The Spark parallelism increase

produced almost no change, which suggests Spark was already using available resources efficiently for this workload. This is proven from the cpu usage charts as well as it Spark is utilizing all the cores available.

Task 2: Real-Time IoT Sensor Data Analytics using Apache Kafka

Part 1: Setup and Environment Configuration

Source : https://github.com/lolitha-lakshan-1/bigdata-assignment/tree/main/Task_2/part-1

1. Installed and configured a Kafka cluster with Docker. The cluster is configured with "kafka", "kafka-exporter", "kafka-connect", "prometheus", "grafana" and "zookeeper" applications.
2. A side cart pod was configured to create the Kafka topics "traffic.raw" and "traffic.processed". The side cart pod will also wait until the Kafka pod get initialized, by monitoring the health check endpoint to prevent premature termination of Kafka Connect.
3. Kafka Connect will support ingesting data through the api and export data to an another medium / storage etc.
4. Kafka Exporter, Prometheus and Grafana will help with the visualization of stats related to the cluster.
5. Following Grafana dashboard panels was developed to cover the aspects of health and performance of the cluster, and the Prometheus was configured as a data source of the dashboard. Grafana has following dashboard panels.

- "Incoming Message Rate" which displays the throughput. A high rate on traffic.raw with low on traffic.processed indicates a processing bottleneck.

```
"sum by (topic) (rate(kafka_topic_partition_current_offset{topic=~\"traffic.*\"}[1m]))"
```

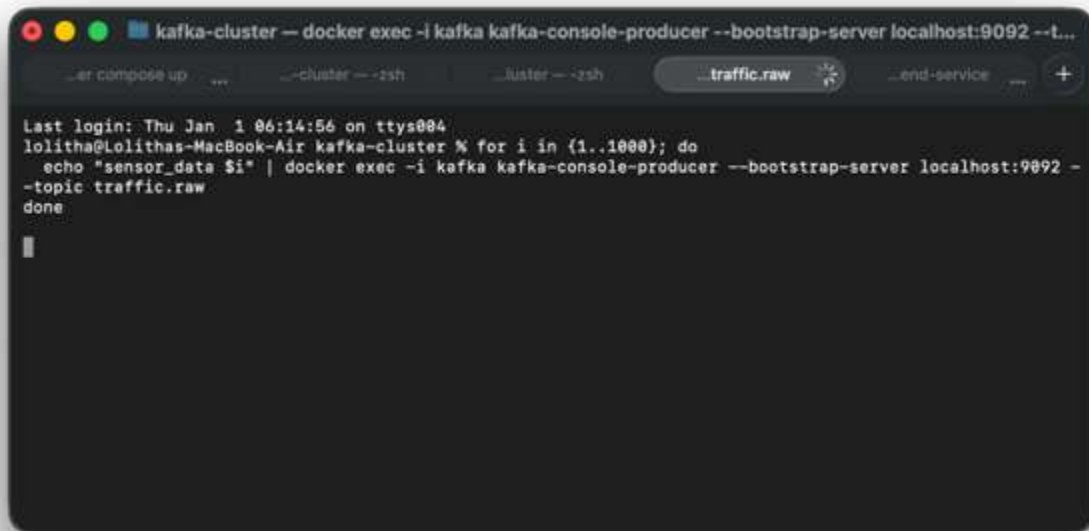
- Total Messages Processed which indicates volume.

```
"sum by (topic) (kafka_topic_partition_current_offset{topic=~\"traffic.*\"})"
```

- Consumer Lag panel indicates how far behind consumers are in reading messages from topics which indicates Processing Health.

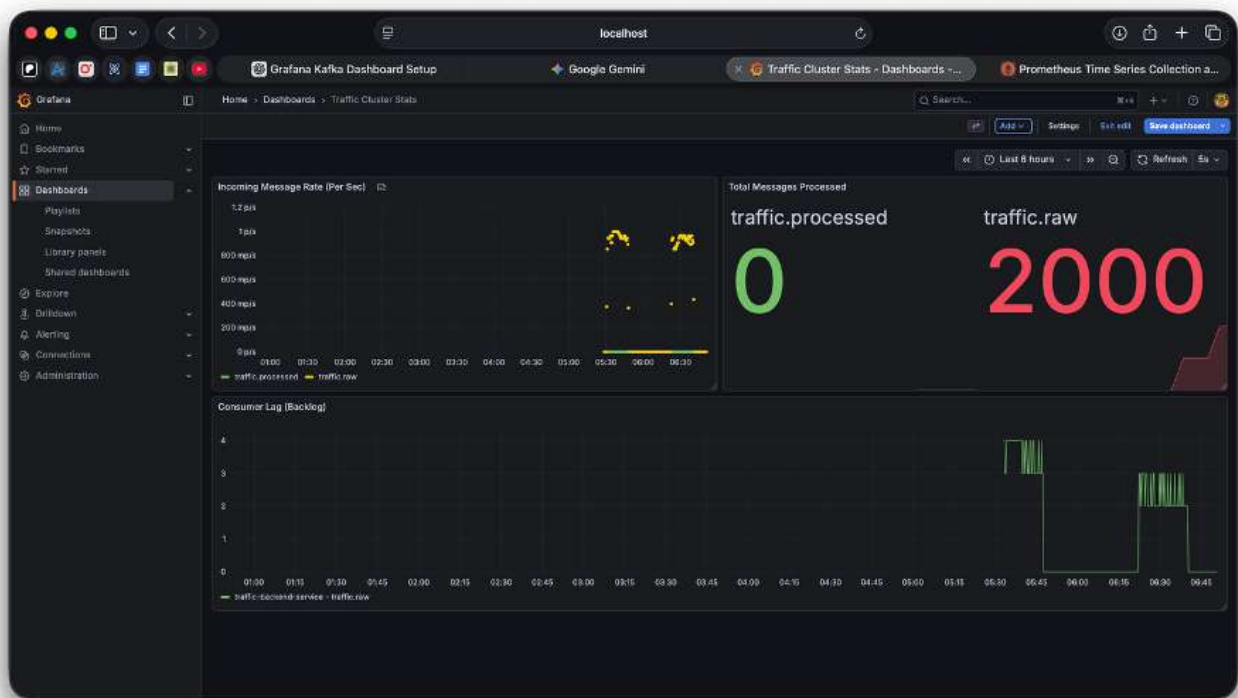
```
"sum by (consumergroup, topic) (kafka_consumergroup_lag{topic=~\"traffic.*\"})"
```

6. Data ingestion was simulated using the command to ingest 1000 rows through a loop.

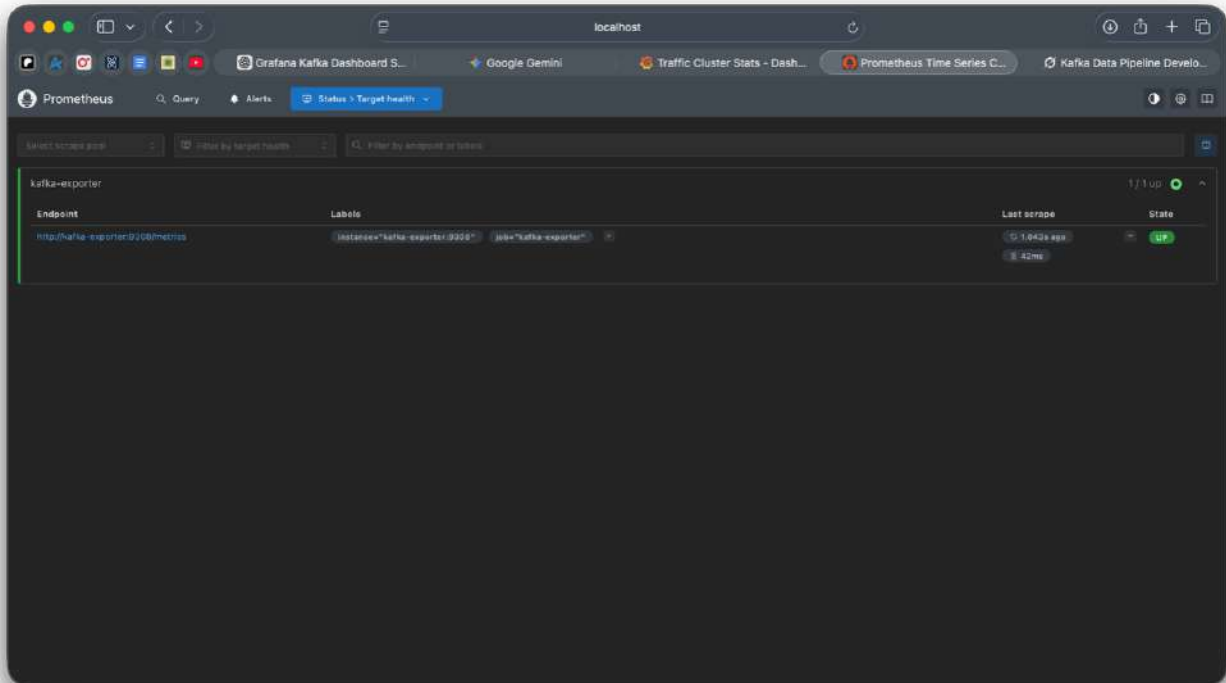


```
kafka-cluster — docker exec -i kafka kafka-console-producer --bootstrap-server localhost:9092 --t...
...er compose up ...
...cluster — -zsh ...cluster — -zsh ...traffic.raw ...end-service +
Last login: Thu Jan 1 06:14:56 on ttys004
lolitha@Lolithas-MacBook-Air kafka-cluster % for i in {1..1000}; do
  echo "sensor_data $i" | docker exec -i kafka kafka-console-producer --bootstrap-server localhost:9092 -
  -topic traffic.raw
done
```

Grafana Dashboard



Prometheus Exporter



Environment Setup

```
kafka-cluster -- zsh -- 105x20

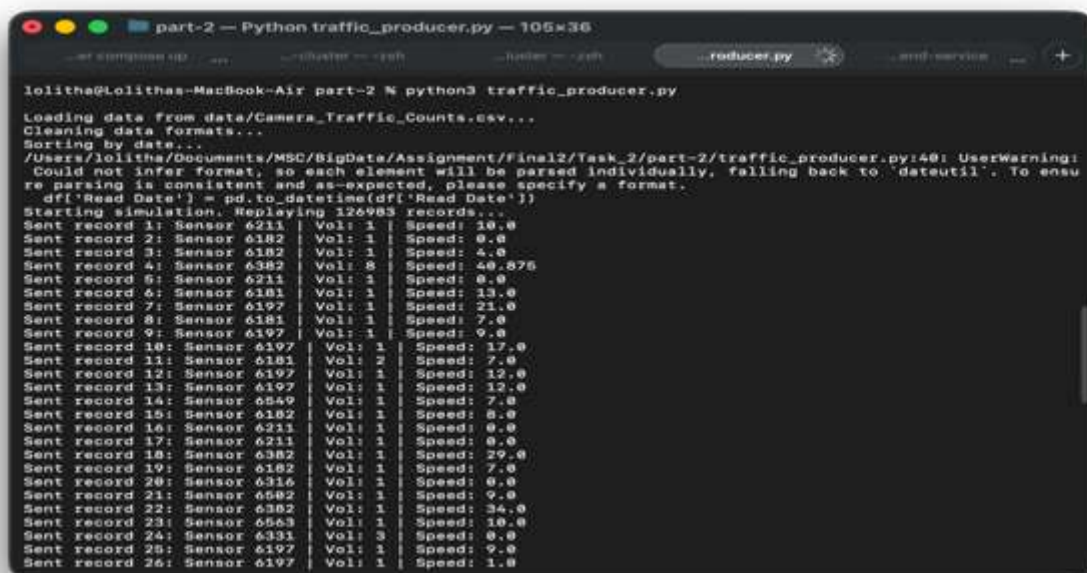
.../part-1 -- zsh  ... er compose up  ...luster -- zsh  ...luster -- zsh  ...end-service  +

Last login: Thu Jan  1 05:32:43 on ttys003
lolitha@Lolithas-MacBook-Air kafka-cluster % docker ps
CONTAINER ID   IMAGE                                COMMAND                  NAMES                  CREATED        STATUS
af7949e140f3   danielqsj/kafka-exporter            "/bin/kafka_exporter..."  kafka-exporter         47 minutes ago Up 47 minutes
095182ee3d67   confluentinc/cp-kafka:7.5.0         "/etc/confluent/dock..."  kafka                  47 minutes ago Up 47 minutes
(healthy)     0.0.0.0:9092->9092/tcp, [::]:9092->9092/tcp
e3e6a74d1d0b   confluentinc/cp-zookeeper:7.5.0     "/etc/confluent/dock..."  zookeeper              47 minutes ago Up 47 minutes
2181/tcp, 2888/tcp, 3888/tcp
75ebae55ef9e   prom/prometheus                     "/bin/prometheus --c..."  prometheus             48 minutes ago Up 47 minutes
0.0.0.0:9090->9090/tcp, [::]:9090->9090/tcp
45221b7f6d39   grafana/grafana                     "/run.sh"                grafana                 48 minutes ago Up 47 minutes
0.0.0.0:3000->3000/tcp, [::]:3000->3000/tcp
1ddd8145fb7f   mongo:latest                         "docker-entrypoint.s..."  mongodb                 3 weeks ago    Up 7 hours
0.0.0.0:27017->27017/tcp, [::]:27017->27017/tcp
lolitha@Lolithas-MacBook-Air kafka-cluster %
```

Part 2: Data Source and Preprocessing

A Python script was written as the Kafka Producer which simulate traffic data every 5 seconds by reading the CSV. This script uses the kafka library and data is sent as JSON messages to the "traffic.raw" kafka topic. Prior to the conversion the script cleans the data by removing commas from the columns (e.g., "1,234" → 1234), also cleans and converts numeric values to integer. Furthermore, drop any rows if column values are missing.

Source: https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_2/part-2/traffic_producer.py



```
lolitha@lolithas-MacBook-Air part-2 % python3 traffic_producer.py
Loading data from data/Camera_Traffic_Counts.csv...
Cleaning data formats...
Sorting by date...
/Users/lolitha/Documents/MSD/BigData/Assignment/Final2/Task_2/part-2/traffic_producer.py:40: UserWarning:
Could not infer format, so each element will be parsed individually, falling back to 'dateutil'. To ensu
re parsing is consistent and as-expected, please specify a format.
  dff['Read Date'] = pd.to_datetime(dff['Read Date'])
Starting simulation. Replaying 126983 records...
Sent record 1: Sensor 6211 | Vol: 1 | Speed: 10.0
Sent record 2: Sensor 6182 | Vol: 1 | Speed: 0.0
Sent record 3: Sensor 6182 | Vol: 1 | Speed: 4.0
Sent record 4: Sensor 6382 | Vol: 8 | Speed: 40.876
Sent record 5: Sensor 6211 | Vol: 1 | Speed: 0.0
Sent record 6: Sensor 6181 | Vol: 1 | Speed: 13.0
Sent record 7: Sensor 6197 | Vol: 1 | Speed: 21.0
Sent record 8: Sensor 6181 | Vol: 1 | Speed: 7.0
Sent record 9: Sensor 6197 | Vol: 1 | Speed: 9.0
Sent record 10: Sensor 6197 | Vol: 1 | Speed: 17.0
Sent record 11: Sensor 6181 | Vol: 2 | Speed: 7.0
Sent record 12: Sensor 6197 | Vol: 1 | Speed: 12.0
Sent record 13: Sensor 6197 | Vol: 1 | Speed: 12.0
Sent record 14: Sensor 6549 | Vol: 1 | Speed: 7.0
Sent record 15: Sensor 6182 | Vol: 1 | Speed: 8.0
Sent record 16: Sensor 6211 | Vol: 1 | Speed: 0.0
Sent record 17: Sensor 6211 | Vol: 1 | Speed: 0.0
Sent record 18: Sensor 6382 | Vol: 1 | Speed: 29.0
Sent record 19: Sensor 6182 | Vol: 1 | Speed: 7.0
Sent record 20: Sensor 6316 | Vol: 1 | Speed: 0.0
Sent record 21: Sensor 6582 | Vol: 1 | Speed: 9.0
Sent record 22: Sensor 6382 | Vol: 1 | Speed: 34.0
Sent record 23: Sensor 6563 | Vol: 1 | Speed: 10.0
Sent record 24: Sensor 6331 | Vol: 3 | Speed: 0.0
Sent record 25: Sensor 6197 | Vol: 1 | Speed: 9.0
Sent record 26: Sensor 6197 | Vol: 1 | Speed: 1.0
```

Part 3: Streaming Data Processing and Analysis

3.1 Data Ingestion

Modified the previously written “traffic_producer.py” script to send data continuously without a delay.

3.2 Processing and Computation

Wrote a real time stream processor using python to consume data from a Kafka topic (traffic.raw), it aggregates data into hourly averages and daily summaries. Thereafter it writes to a PostgreSQL database and publish the processed data to the "traffic.processed" Kafka topic.

Source : https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_2/part-3/traffic_processor.py

The "traffic_processor.py" stream processor has following main helper functions to calculate the statistics.

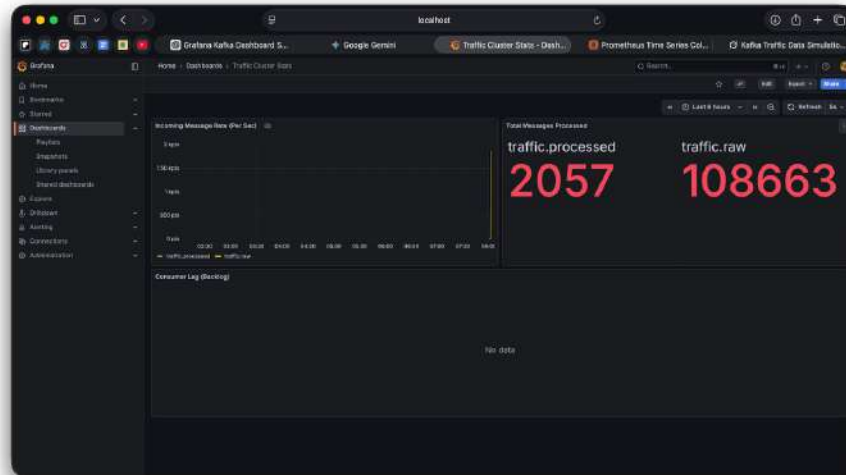
The `update_running_totals` function processes each incoming traffic record. For the given record, it adds the record's volume to `sum_vol` and the speed to `sum_speed` within the `hourly_stats` dictionary entry for that sensor. It also increments the count of records seen for that sensor in the current hour. Simultaneously, it compares the current volume against the `peak_volume` in `daily_state`. If higher, it updates the `peak_volume` and sets `peak_sensor` to the current `sensor_id`. Finally, it adds the `sensor_id` to both the `active_sensors_today` set (tracking sensors reporting on the current day) and the `all_known_sensors` set (tracking every sensor ever encountered).

The `calculate_hourly_averages` function is called for each new hour. It iterates over the accumulated `hourly_stats` dictionary from the current hour and computes `avg_vol` as `sum_vol / count` and `avg_speed` as `sum_speed / count` for each sensor with data. It builds a list of `db_rows` for database insertion and creates corresponding `kafka_messages`. After processing, `hourly_stats` is reset to an empty dictionary for the new hour.

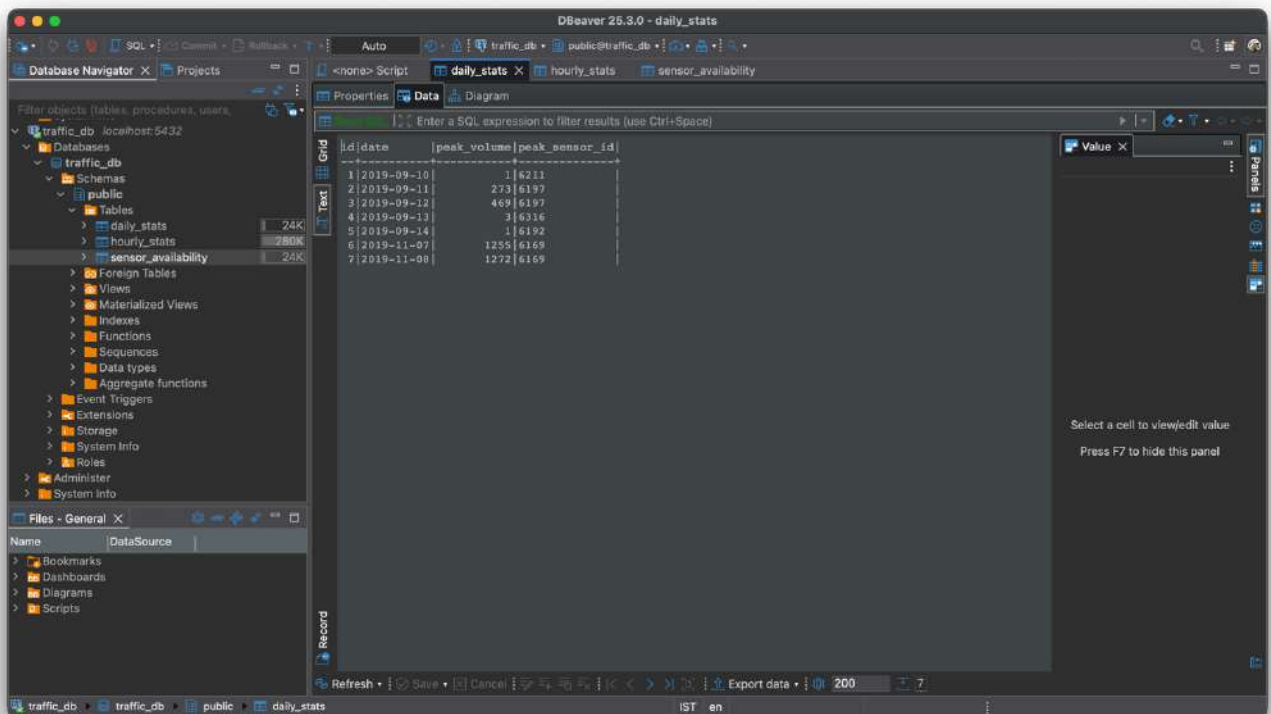
The `calculate_daily_summary` function executes when the date changes. It uses the data accumulated in `daily_state` to determine the day's `peak_volume` and associated `peak_sensor`, counts the number of `active_sensors_today`, and computes `availability_pct` as the percentage of `active_sensors_today` relative to the total size of `all_known_sensors`. It returns a summary dictionary containing the date, peak values, active count, and availability percentage. Subsequently, the daily tracking fields—`peak_volume`, `peak_sensor`, and `active_sensors_today`—are reset, while `all_known_sensors` persists across days.

3.3 Output and Persistence

Kafka Data Processed Topic



Data Stored in Postgres DB



DBeaver 25.3.0 - hourly_stats

Filter objects (tables, procedures, users, ...)

traffic_db - localhost:5432

- traffic_db
 - Schemas
 - public
 - Tables
 - daily_stats (24K)
 - hourly_stats (280K)
 - sensor_availability (24K)
 - Foreign Tables
 - Views
 - Materialized Views
 - Indexes
 - Functions
 - Sequences
 - Data types
 - Aggregate functions
 - Event Triggers
 - Extensions
 - Storage
 - System Info
 - Roles
 - Administrator
 - System Info

Files - General

Name | DataSource

Bookmarks

Dashboards

Diagrams

Scripts

Grid

| id | sensor_id | window_start | avg_volume | avg_speed | record_count |
|----|-----------|-------------------------|------------|-----------|--------------|
| 1 | 6211 | 2019-09-10 00:00:00.000 | 1.0 | 10.0 | 1 |
| 2 | 6182 | 2019-09-10 01:00:00.000 | 1.0 | 0.0 | 1 |
| 3 | 6182 | 2019-09-10 02:00:00.000 | 1.0 | 4.0 | 1 |
| 4 | 6211 | 2019-09-10 05:00:00.000 | 1.0 | 0.0 | 1 |
| 5 | 6197 | 2019-09-10 09:00:00.000 | 1.0 | 21.0 | 1 |
| 6 | 6197 | 2019-09-10 12:00:00.000 | 1.0 | 13.0 | 2 |
| 7 | 6197 | 2019-09-10 15:00:00.000 | 1.0 | 12.0 | 1 |
| 8 | 6197 | 2019-09-10 16:00:00.000 | 1.0 | 12.0 | 1 |
| 9 | 6182 | 2019-09-10 18:00:00.000 | 1.0 | 8.0 | 1 |
| 10 | 6211 | 2019-09-10 19:00:00.000 | 1.0 | 0.0 | 1 |
| 11 | 6211 | 2019-09-10 20:00:00.000 | 1.0 | 0.0 | 1 |
| 12 | 6182 | 2019-09-11 00:00:00.000 | 1.0 | 7.0 | 1 |
| 13 | 6316 | 2019-09-11 00:00:00.000 | 1.0 | 0.0 | 1 |
| 14 | 6502 | 2019-09-11 00:00:00.000 | 1.0 | 9.0 | 1 |
| 15 | 6331 | 2019-09-11 07:00:00.000 | 3.0 | 0.0 | 1 |
| 16 | 6197 | 2019-09-11 08:00:00.000 | 1.0 | 5.0 | 2 |
| 17 | 6211 | 2019-09-11 10:00:00.000 | 1.0 | 9.0 | 1 |
| 18 | 6331 | 2019-09-11 10:00:00.000 | 1.0 | 6.0 | 1 |
| 19 | 6197 | 2019-09-11 16:00:00.000 | 87.0 | 34.552 | 1 |
| 20 | 6331 | 2019-09-11 11:00:00.000 | 1.0 | 31.0 | 1 |
| 21 | 6331 | 2019-09-11 12:00:00.000 | 1.0 | 18.0 | 1 |
| 22 | 6197 | 2019-09-11 13:00:00.000 | 159.0 | 31.676 | 2 |
| 23 | 6197 | 2019-09-11 16:00:00.000 | 1.0 | 7.0 | 1 |
| 24 | 6331 | 2019-09-11 20:00:00.000 | 1.0 | 8.0 | 1 |
| 25 | 6211 | 2019-09-11 20:00:00.000 | 1.0 | 10.0 | 1 |
| 26 | 6331 | 2019-09-11 21:00:00.000 | 3.0 | 11.333 | 1 |
| 27 | 6140 | 2019-09-11 22:00:00.000 | 1.0 | 10.0 | 1 |
| 28 | 6197 | 2019-09-12 00:00:00.000 | 1.0 | 12.0 | 1 |
| 29 | 6546 | 2019-09-12 05:00:00.000 | 1.0 | 10.0 | 1 |
| 30 | 6331 | 2019-09-12 06:00:00.000 | 15.0 | 13.6 | 1 |
| 31 | 6197 | 2019-09-12 07:00:00.000 | 1.0 | 13.0 | 1 |
| 32 | 6197 | 2019-09-12 08:00:00.000 | 469.0 | 36.74 | 1 |
| 33 | 6331 | 2019-09-12 09:00:00.000 | 42.0 | 14.571 | 1 |
| 34 | 6502 | 2019-09-12 09:00:00.000 | 1.0 | 6.0 | 1 |
| 35 | 6211 | 2019-09-12 10:00:00.000 | 1.0 | 5.0 | 1 |
| 36 | 6331 | 2019-09-12 11:00:00.000 | 30.0 | 18.6 | 1 |

Refresh | Save | Cancel | Export data | 200 | 200+

DBeaver 25.3.0 - sensor_availability

Filter objects (tables, procedures, users, ...)

traffic_db - localhost:5432

- traffic_db
 - Schemas
 - public
 - Tables
 - daily_stats (24K)
 - hourly_stats (280K)
 - sensor_availability (24K)
 - Foreign Tables
 - Views
 - Materialized Views
 - Indexes
 - Functions
 - Sequences
 - Data types
 - Aggregate functions
 - Event Triggers
 - Extensions
 - Storage
 - System Info
 - Roles
 - Administrator
 - System Info

Files - General

Name | DataSource

Bookmarks

Dashboards

Diagrams

Scripts

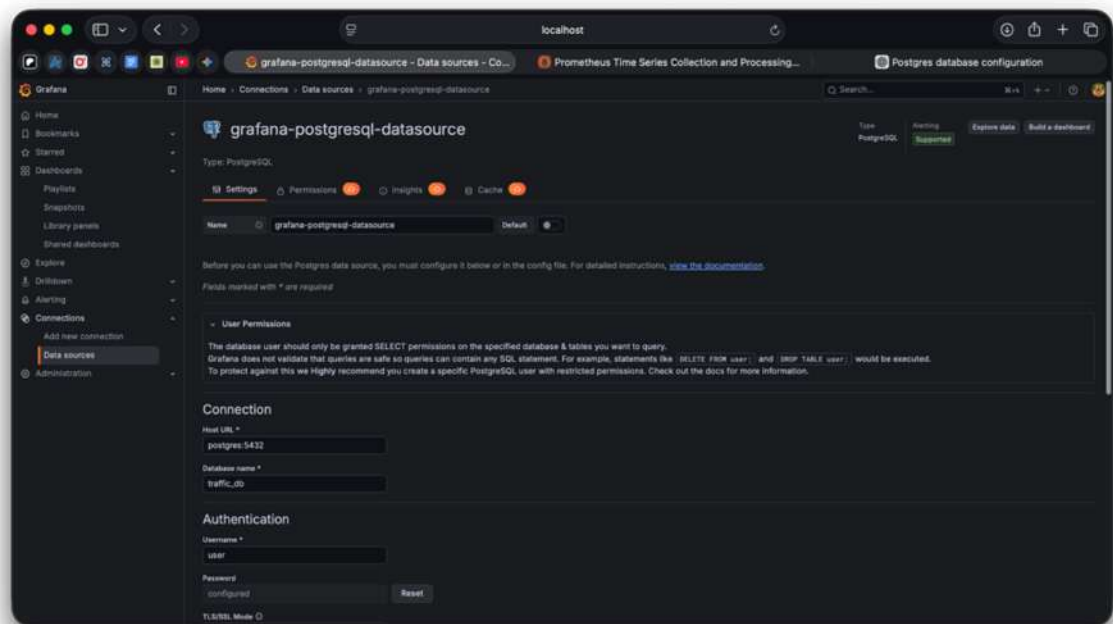
Grid

| id | date | total_sensors_seen | availability_percent |
|----|------------|--------------------|----------------------|
| 1 | 2019-09-10 | 3 | 100.0 |
| 2 | 2019-09-11 | 7 | 100.0 |
| 3 | 2019-09-12 | 6 | 75.0 |
| 4 | 2019-09-13 | 8 | 72.72727272727273 |
| 5 | 2019-09-14 | 2 | 18.181818181818183 |
| 6 | 2019-11-07 | 35 | 94.5945945945946 |
| 7 | 2019-11-08 | 41 | 93.18181818181817 |

Refresh | Save | Cancel | Export data | 200 | 7

Part 4: Visualization and Reporting

1. First updated the Grafana datasource with the Postgres datasource and added the db name and db credentials.

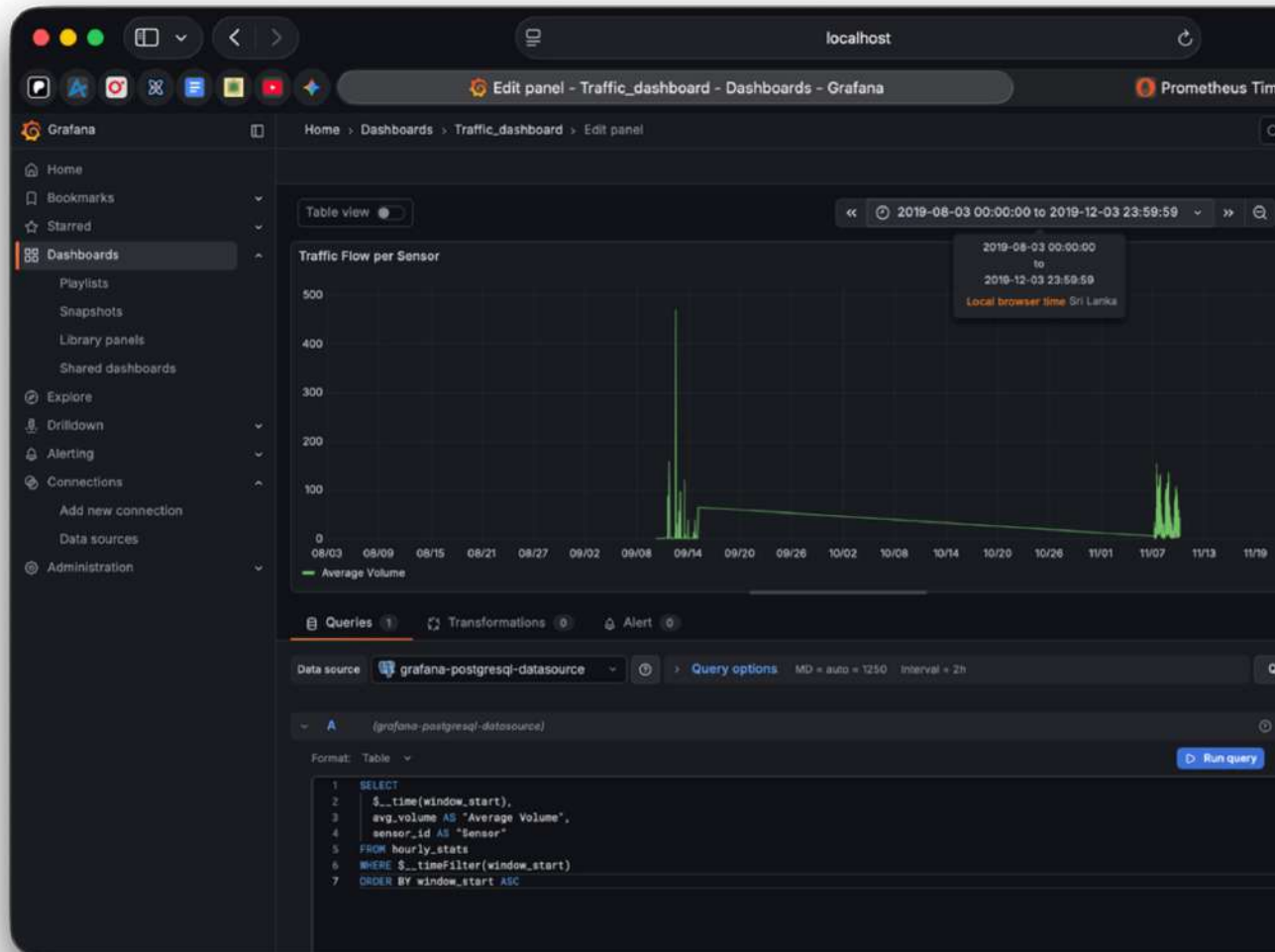


2. Created the dashboard panels using the UI and following queries.



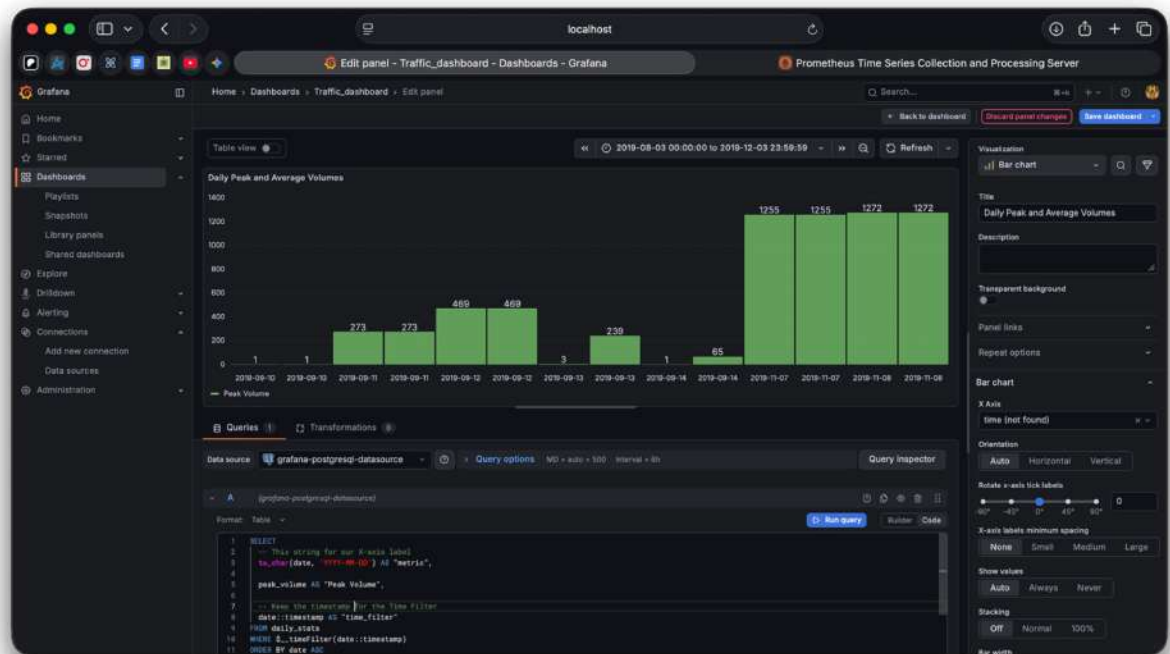
Source: https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_2/part-4/kafka-cluster/grafana/provisioning/dashboards/traffic_dashboard.json

Time-series chart for traffic flow per sensor



This query works by selecting the `window_start` as the timestamp and the `avg_volume` as the plotted value. The `sensor_id` is cast to text using, which forces Grafana to treat the ID as a label rather than a number, resulting in a distinct line for every sensor.

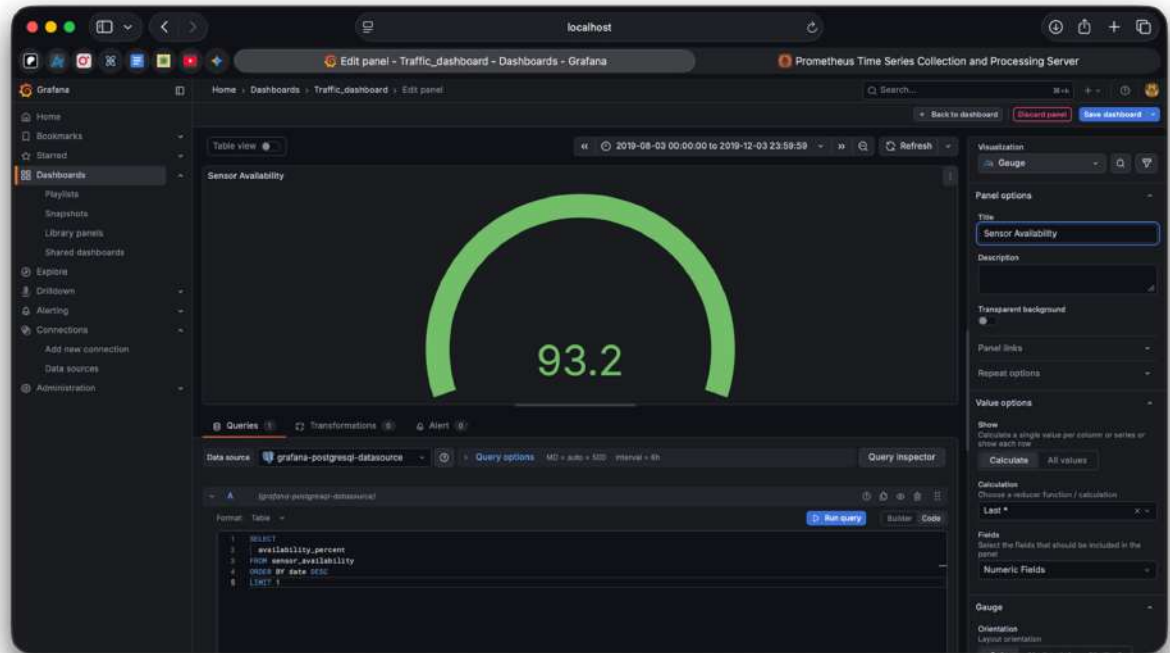
Bar chart for daily peak and average volumes



```
SELECT
to_char(date, 'YYYY-MM-DD') AS "metric",
peak_volume AS "Peak Volume",
date::timestamp AS "time_filter"
FROM daily_stats
WHERE $__timeFilter(date::timestamp)
ORDER BY date ASC
```

Since Grafana's time-axis requires precise timestamps rather than simple dates, date column had to be cased to place the bars correctly on the X-axis.

Gauge indicators or tables for sensor availability and Health



```
SELECT
  availability_percent
FROM sensor_availability
ORDER BY date DESC
LIMIT 1
```

The Sensor Availability gauge was created with the above query which will fetch most recent record by ordering the data descending.

Task 3: Graph Databases using Neo4j and Docker Objective

Part 1 — Understanding Graph Databases

Fraud Detection In Financial Services

Major financial institutions like JPMorgan Chase utilize graph databases to uncover sophisticated fraud rings or money laundering schemes. The traditional relational databases struggle in this context as indirect connections would require complex JOIN queries to uncover these sorts of data. In a graph data base Nodes will represent entities such as Customers, Accounts, ATMs as well as meta data such as ip addresses, device ids and phone numbers.

Graph relationships will include actions such as `TRANSFERRED_TO`, `LOGGED_IN_WITH`, or `SHARES_DEVICE`. Therefore, a graph database will help investigators to instantly traverse through the network to spot non-obvious patterns, such as 20 seemingly unrelated accounts all accessing banking services from the same mobile device ID. In a graph database these kinds of patterns would be easy to spot and discover than a traditional relational database. In a traditional database this information will be buried under columns, rows and may be distributed across tables.

- Useful Query: Cycle Detection like find a path where money leaves Account A, passes through 5 intermediaries, and returns to an account held by A's known associate.

Source: [Neo4j Financial Services Use Cases](#) & [JPMorgan Chase AI Research](#)

Food Discovery

Uber Eats built a Food Knowledge Graph to address the zero-result problem which occurs when a specific search term yields no exact match, which could user to abandon the app. A standard database would simply return an empty results for non-existing search term, but a graph structure has the capability to understand the *semantic* relationships between food items, allowing for intelligent substitutions.

In a graph database, nodes can be utilized save entities such as menu items, cuisines, food attributes such as Spicy, Vegan. Furthermore, relationships can be utilized to define associations of the nodes such as `IS_A_TYPE_OF` or `OFTEN_ORDERED_WITH`.

With these relationships graph databases will allow to expand the query to look for associated food items if no direct food items were found. For an example If a user searches for a specific dish that isn't available, the graph knows that "Udon" is related to "Ramen" via the "Japanese Noodle" parent node, allowing the app to recommend a relevant alternative instantly.

- Useful Query: Semantic Similarity (e.g., Find all restaurants serving items connected to the Japanese Cuisine)

Source: [Uber Engineering Blog: Food Discovery with Uber Eats](#)

Contextual Search & Recommendations In Ecommerce

eBay developed a distributed knowledge graph (internally called "Akutan") to better understand the massive inventory of over 1 billion items, moving beyond simple keyword matching. Relational databases are rigid and can struggle to define with the highly variable attributes of diverse products (e.g., a "screen size" matters for phones but not for sneakers), whereas a graph flexibly links products to real-world concepts. In an ecommerce site

Products (iPhone 12), Events (Super Bowl), and Entities (Celebrities, Brands) can be stored as graph database nodes and connections such as like `COMPATIBLE_WITH`, `USED_BY`, or `RELEASED_DURING` can be expressed using graph database relationships. These relationships would help with complex questions where keyword search will fail for an example a specific phone case `IS_COMPATIBLE_WITH` a specific phone model, eBay can filter out thousands of irrelevant results that just happen to have similar keywords in the title.

- Useful Query: Pathfinding (e.g., Return all distinct accessories linked via `COMPATIBLE_WITH` edges to the specific electronics item the user is currently viewing).

Source: [eBay Innovation Blog: Akutan Knowledge Graph](#)

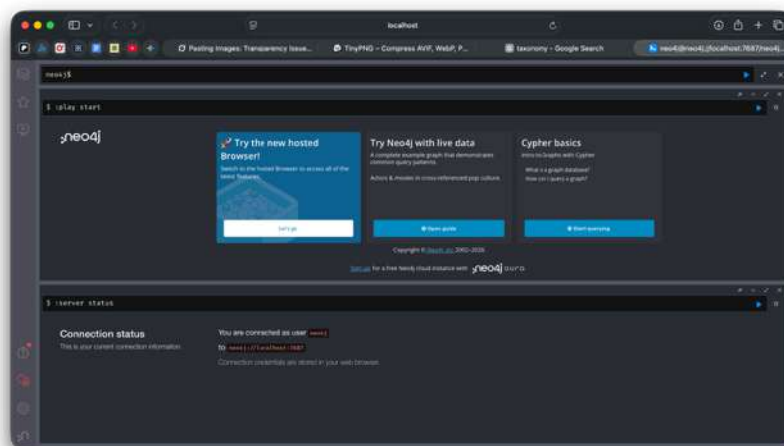
Part 2 — Implementation with Neo4j (Docker Setup)

Implementation

```
Task_3 — -zsh — 80x21

...lgData/Assignment/Final2/Task_3 — -zsh ... ...Assignment/Final2/Task_3 — -zsh +

lolitha@Lolithas-MacBook-Air Task_3 % docker run -d \
  --name "neo4j-graphdb" \
  -p7474:7474 -p7687:7687 \
  -e NEO4J_AUTH=neo4j/test12345 \
  [ "neo4j:latest"
c268b47b602e5a7c002ef38c5f39557cee8df00cf06c2fd3e07da8fdd126abab
lolitha@Lolithas-MacBook-Air Task_3 %
```



Dataset

The dataset simulates a financial fraud ring. It connects people, accounts, and digital footprint to reveal suspicious activity.

Following nodes are present in the dataset,

- Customers: 5 individuals (C1–C5).
- Accounts: 6 bank accounts (A1–A6).
- Infrastructure: 3 Devices (Phones/Tablets) and 2 IP Addresses.

And Nodes are connected using following relationships,

- Customers OWNS accounts.
- Customers LOGGED_IN_WITH devices and LOGGED_IN_FROM IP addresses.
- Accounts TRANSFERRED_TO other accounts with amount and currency.

As an overview the C1 Customer controls two accounts, A1 and A2, and money transfers to A1 and moves step by step through other customers' accounts, then returns to A1, which creates a clear circular transfer loop. The login links make the connections feel intentional because C1 and C2 share the same device and the same IP address, C3 and C4 share another device and IP, and C5 also logs in from the same IP as C1 and C2 which suggests a money mule.

```

fraud_data_Cnames_customers_login.cypher
MATCH (n)
DETACH DELETE n;

CREATE
(c1:Customer {id: 'C1', name: 'C1'})-[:OWNS]->(a1:Account {id: 'A1', accountNo: 'ACC-001'}),
(c1)-[:OWNS]->(a2:Account {id: 'A2', accountNo: 'ACC-002'}),

(c2:Customer {id: 'C2', name: 'C2'})-[:OWNS]->(a3:Account {id: 'A3', accountNo: 'ACC-003'}),
(c3:Customer {id: 'C3', name: 'C3'})-[:OWNS]->(a4:Account {id: 'A4', accountNo: 'ACC-004'}),
(c4:Customer {id: 'C4', name: 'C4'})-[:OWNS]->(a5:Account {id: 'A5', accountNo: 'ACC-005'}),
(c5:Customer {id: 'C5', name: 'C5'})-[:OWNS]->(a6:Account {id: 'A6', accountNo: 'ACC-006'}),

(d1:Device {id: 'D1', deviceId: 'dev-xyz', type: 'Phone'}),
(d2:Device {id: 'D2', deviceId: 'dev-abc', type: 'Phone'}),
(d3:Device {id: 'D3', deviceId: 'dev-789', type: 'Tablet'}),

(ip1:IPAddress {id: 'IP1', address: '10.0.0.10'}),
(ip2:IPAddress {id: 'IP2', address: '10.0.0.20'}),

// Money transfers (fraud ring pattern)
(a1)-[:TRANSFERRED_TO {amount: 2500, currency: 'USD'}]->(a3),
(a3)-[:TRANSFERRED_TO {amount: 2600, currency: 'USD'}]->(a4),
(a4)-[:TRANSFERRED_TO {amount: 2700, currency: 'USD'}]->(a5),
(a5)-[:TRANSFERRED_TO {amount: 2800, currency: 'USD'}]->(a6),
(a2)-[:TRANSFERRED_TO {amount: 5000, currency: 'USD'}]->(a6),
(a6)-[:TRANSFERRED_TO {amount: 2400, currency: 'USD'}]->(a1),

(c1)-[:LOGGED_IN_WITH]->(d1),
(c2)-[:LOGGED_IN_WITH]->(d1),
(c3)-[:LOGGED_IN_WITH]->(d2),
(c4)-[:LOGGED_IN_WITH]->(d2),
(c5)-[:LOGGED_IN_WITH]->(d3),

(c1)-[:LOGGED_IN_FROM]->(ip1),
(c2)-[:LOGGED_IN_FROM]->(ip1),
(c5)-[:LOGGED_IN_FROM]->(ip1),
(c3)-[:LOGGED_IN_FROM]->(ip2),
(c4)-[:LOGGED_IN_FROM]->(ip2);

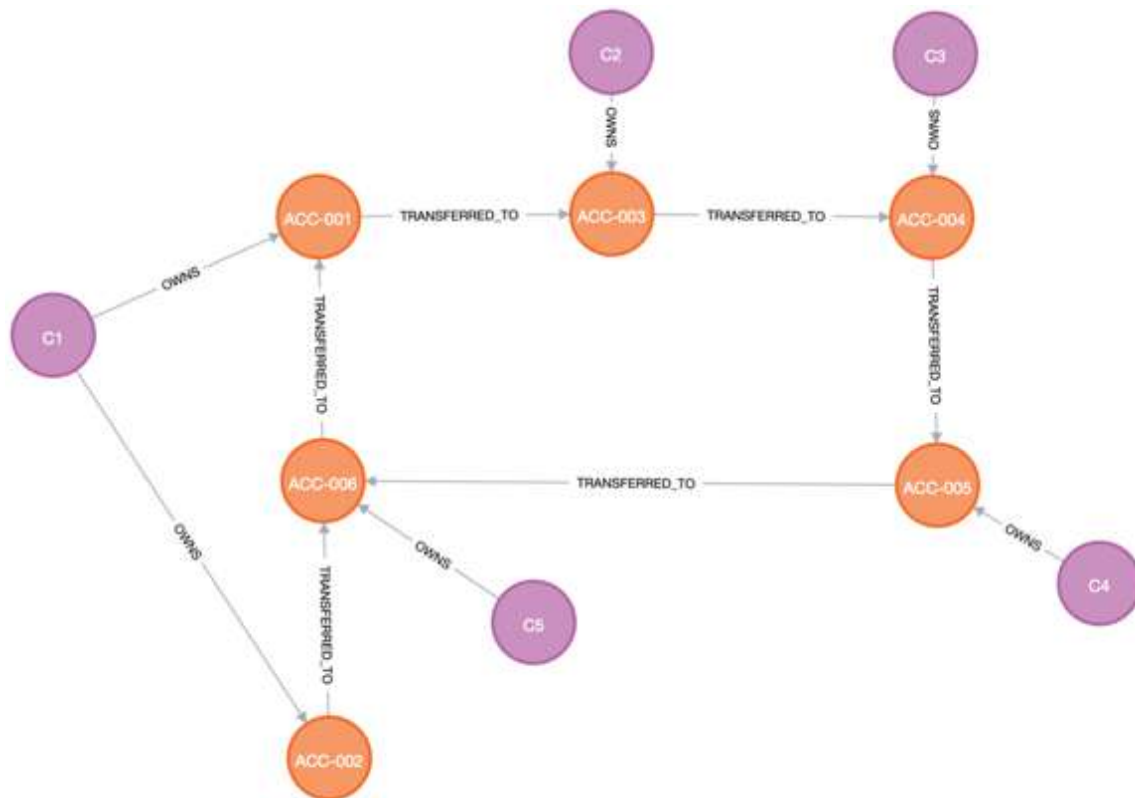
```

Queries

1. Find Associates of C1Customer

```
MATCH (startCustomer:Customer {name: 'C1'})-[:OWNS]->(startAccount:Account)
MATCH paths = (startAccount)-[:TRANSFERRED_TO*1..4]-(:Account)<-[:OWNS]-(otherCustomer:Customer)
RETURN paths, otherCustomer;
```

The query is trying to answer “Which other customers are connected to C1 through fund transfers?” and it does this by starting from C1’s accounts and following 1– 4 TRANSFERRED_TO hops between accounts to find any other customers who receive money directly or indirectly.



The result shows that money flowing out of ACC-001 and ACC-002 and eventually reaches accounts owned by C2, C3, C4, and C5. This means all four of those customers are transactionally connected to C1, forming a small network of potentially related or suspicious activity.

2. Most connected node

```
MATCH (node) RETURN labels(node) AS nodeLabels, node.id AS nodeId, COUNT { (node)--() } AS degree ORDER BY degree DESC LIMIT 5;
```

This query tries to find which nodes act as hubs in the fraud graph by counting how many relationships each node has. It counts all connections and then returns the top 5 most connected nodes.

| nodeLabels | nodeId | degree |
|--------------|--------|--------|
| ["Account"] | "A6" | 4 |
| ["Customer"] | "C1" | 4 |
| ["Customer"] | "C2" | 3 |
| ["Account"] | "A3" | 3 |
| ["Account"] | "A1" | 3 |

The result shows that account A6 and customer C1 both have degree of 4, meaning they each touch four other nodes, making them strong candidates as suspicious entities. Accounts A3 and A1, and customer C2, also have relatively high degrees (3), indicating they participate in multiple transfers or relationships.

- A6 – owned by customer C5; receives money from A5 and A2 and sends money to A1.
- C1 – owns accounts A1 and A2, has LOGGED_IN_WITH → D1 and has LOGGED_IN_FROM → IP1.
- C2 – owns account A3, has LOGGED_IN_WITH → D1 and has LOGGED_IN_FROM → IP1.
- A3 – owned by C2, receives money from A1, sends money to A4.
- A1 – owned by C1, receives money from A6, sends money to A3.

3. Customers sharing both a device and an IP

This query looks for pairs of customers who share both the same device and the same IP address, which is a strong signal that they might be the same person or closely colluding accounts.

```
MATCH
  (customer1:Customer)-[:LOGGED_IN_WITH]->(device:Device)<-[:LOGGED_IN_WITH]-
(customer2:Customer),
  (customer1)-[:LOGGED_IN_FROM]->(ip:IPAddress)<-[:LOGGED_IN_FROM]-(customer2)
WHERE customer1 <> customer2
RETURN DISTINCT
  customer1.name,
  customer2.name,
  device.deviceId,
  ip.address
ORDER BY customer1.name, customer2.name;
```

| customer1.name | customer2.name | device.deviceId | ip.address |
|----------------|----------------|-----------------|-------------|
| "C1" | "C2" | "dev-xyz" | "10.0.0.10" |
| "C2" | "C1" | "dev-xyz" | "10.0.0.10" |
| "C3" | "C4" | "dev-abc" | "10.0.0.20" |
| "C4" | "C3" | "dev-abc" | "10.0.0.20" |

The result shows C1–C2 sharing dev-xyz and 10.0.0.10, and C3–C4 sharing dev-abc and 10.0.0.20. These pairs should get high interest from a fraud perspective because multiple customer profiles appear to be operated from the same physical device and network.

Screenshots

The screenshot shows the Neo4j Browser interface. On the left, the 'Database Information' sidebar is visible, showing the database name 'neo4j', node labels ('(ID)', 'Account', 'Customer', 'Device', 'IPAddress'), relationship types ('LOGGED_IN_FROM', 'LOGGED_IN_WITH', 'OWNED', 'TRANSFERRED_TO'), and property keys ('accountNo', 'address', 'amount', 'currency', 'deviceId', 'id', 'name', 'type'). The 'Connected as' section shows the username 'neo4j' and roles. The 'DBMS' section shows the cluster role 'primary' and version '2020.10.1'.

The main query editor shows the following Cypher query:

```
neo4j$ MATCH (node) RETURN labels(node) AS nodeLabels, node.id AS nodeId, COUNT { (node)...
```

The results are displayed in a table with columns: nodeLabels, nodeId, degree.

| nodeLabels | nodeId | degree |
|--------------|--------|--------|
| ["Account"] | "A6" | 4 |
| ["Customer"] | "C1" | 4 |
| ["Customer"] | "C2" | 3 |
| ["Account"] | "A3" | 3 |
| ["Account"] | "A1" | 3 |

The screenshot shows the Neo4j Browser interface. On the left, the 'Database Information' sidebar is visible, showing the database name 'neo4j', node labels ('(ID)', 'Account', 'Customer', 'Device', 'IPAddress'), relationship types ('LOGGED_IN_FROM', 'LOGGED_IN_WITH', 'OWNED', 'TRANSFERRED_TO'), and property keys ('accountNo', 'address', 'amount', 'currency', 'deviceId', 'id', 'name', 'type'). The 'Connected as' section shows the username 'neo4j' and roles. The 'DBMS' section shows the cluster role 'primary' and version '2020.10.1'.

The main query editor shows the following Cypher query:

```
neo4j$ MATCH (node) RETURN labels(node) AS nodeLabels, node.id AS nodeId, COUNT { (node)--() } AS degree ORDER BY degree DESC LIMIT 5;
```

The results are displayed in a table with columns: customer1.name, customer2.name, device.deviceId, ip.address.

| customer1.name | customer2.name | device.deviceId | ip.address |
|----------------|----------------|-----------------|-------------|
| "C1" | "C2" | "dev-xyz" | "10.0.0.10" |
| "C2" | "C1" | "dev-xyz" | "10.0.0.10" |
| "C3" | "C4" | "dev-abc" | "10.0.0.20" |
| "C4" | "C3" | "dev-abc" | "10.0.0.20" |

Task 4: Watermarks in Real-Time Stream Processing using Apache Kafka and Apache Flink (Docker Setup)

Part 1 — Understanding Watermarks

*4.1.1 What are **watermarks** in stream processing?*

The real time systems consume never ending stream of data. This data is used to generate outputs such as sums, averages or counts. This is usually done by grouping the events into time windows. However, due to unforeseen circumstances data can arrive out of order. A watermark is a rule which specifies whether the time window for processing the data elapsed or not, rather than waiting forever for late events to arrive.

For an example imagine a city traffic system where the road sensors send the number of cars every minute. Due to a slow network a reading measured at 10:01 might arrive late at 10:03. Without watermarks, the system will close the 10:00–10:01 window too early which would result missing the late reading thus resulting wrong traffic stats,

With watermarks, we can say specify wait up to 30 seconds for late data, then close the window and publish the result which gives a good balance between freshness and accuracy specially for live dashboards.

4.1.2 Watermark Generation Strategies

The three main strategies for generating watermarks are Periodic, Punctuated, and Event-Time-Based.

Periodic watermarks

The periodic watermark strategy generates watermarks at fixed intervals based on processing time or the system clock (e.g., every 100ms). The system calculates the watermark by finding the maximum event time observed during the period and then subtracting a pre-configured 'allowed lateness' delay. This technique ensures that the system can consistently advance without stalling, even when the incoming data stream is inactive or when data is sparse.

Punctuated watermarks

Punctuated watermarks are created only when a special marker event appears in the data stream. This marker event carries a flag or metadata which tells the stream processor that time has elapsed. Therefore, punctuated watermarks do not use a timer; instead, the stream processor waits for these marker events and uses the marker event metadata to generate the new watermark. However, the system will stop advancing the time (stall) if these marker events do not arrive as expected.

Event-time-based watermarks

Event-Time-Based Watermarks are based on the actual event time in the data record itself. e.g. when a social media post was actually made. Rather than considering event received time. Therefore, system is focused on measuring the actual timeline of events, regardless of network delays or out-of-order arrival. The watermark advances only when new, later-timestamped data arrives.

$$\text{Watermark} = \text{Max Event Time Seen} - \text{Fixed Delay}$$

| Strategy | Pros | Cons |
|------------------------------------|---|---|
| Periodic Watermarks | <ul style="list-style-type: none"> • Low stall risk • Low processing overhead, as checks occur only at set intervals. • Skew can be easily adjusted to balance accuracy and latency. | <ul style="list-style-type: none"> • Always introduces a delay, increasing latency. • Cannot handle extremely late events. |
| Punctuated Watermarks | <ul style="list-style-type: none"> • Provides the fastest possible update. | <ul style="list-style-type: none"> • High Stall Risk and High Overhead. • Requires checking every record for the marker flag. • Source Dependency: Relies on the producer to correctly generate the marker events. |
| Event-Time-Based Watermarks | <ul style="list-style-type: none"> • Essential for consistent results regardless of network delays or processing lag. | <ul style="list-style-type: none"> • Requires handling out-of-order data, which introduces complexity in stream processor logic and state management. |

4.1.3) Handling late or out-of-order events

Periodic Watermarks

Periodic Watermarks closes the time window as soon as duration is elapsed. The time is based on the system clock rather than the event time therefore the lagging events are considered late. If the network lags watermarks will be generated continuously which would result dropping valid data. To mitigate these following strategies can be implemented.

- Tune Allowed Lateness. In some systems it is possible to keep a window state around for a specific time after the watermark has passed. If a late event arrives during this grace period, the window is re-opened, the calculation is updated, and a corrected result is emitted.
- Implement a Dead Letter Queue and process the data later.

Punctuated Watermarks

The biggest risk with punctuated watermarks is that if the data source goes quiet or marker events does not arrive, the watermark will stop moving which would result a stagnant window.

To mitigate this following can be done,

- Implement a timeout, window will move after a timeout occurs.
- If the source sends late event after the mark event handling that with a Dead Letter Queue and process later.
- Use Heartbeats events and move the window when the data stream is empty.

Event-Time-Based Watermarks

For Event-Time-Based Watermarks it is possible to,

- Configure the system to define a specific duration of "allowed lateness" after the watermark has passed the window's end time and handle late date. xw
- Dead Letter Queue and process later.

What is best for best for real-world social media streams?

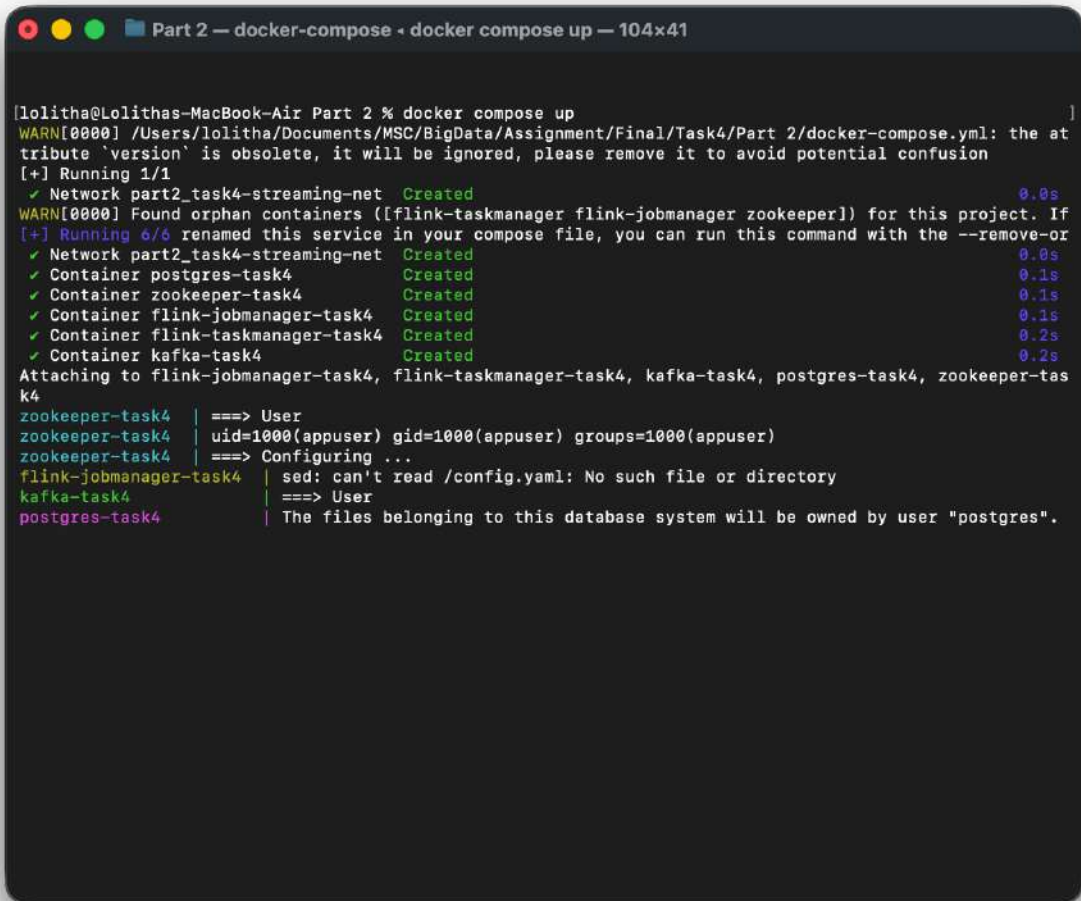
The social media streaming platform have high-volume, spikes of data and globally distributed therefore we should expect to handle bursty and variable network latency of data.

The best strategy to handle this would be Event-Time-Based Watermarks with Bounded Out-of-Orderness.

Since it tracks the true progress of the event stream, ensuring that results are based on the actual time the events occurred, despite events arriving out of order. The "Fixed Delay" parameter will allow to balance the need of accuracy and latency where for social media platforms the accuracy is most probably eventually consistent.

Part 2 — Implementation using Docker (Kafka + Flink)

Step 1 — Environment Setup

A terminal window titled "Part 2 — docker-compose - docker compose up — 104x41" showing the output of the command "docker compose up". The output includes a warning about an obsolete 'version' attribute, a list of created containers and networks with their creation times, and a message about orphan containers. It also shows the configuration of the zookeeper-task4 container, including user and group settings, and a message about the files belonging to the postgres database system.

```
[lolitha@Lolithas-MacBook-Air Part 2 % docker compose up]
WARN[0000] /Users/lolitha/Documents/MS/BigData/Assignment/Final/Task4/Part 2/docker-compose.yml: the at
tribute 'version' is obsolete, it will be ignored, please remove it to avoid potential confusion
[+] Running 1/1
  ✓ Network part2_task4-streaming-net Created 0.0s
WARN[0000] Found orphan containers ([flink-taskmanager flink-jobmanager zookeeper]) for this project. If
[+] Running 6/6 renamed this service in your compose file, you can run this command with the --remove-or
  ✓ Network part2_task4-streaming-net Created 0.0s
  ✓ Container postgres-task4 Created 0.1s
  ✓ Container zookeeper-task4 Created 0.1s
  ✓ Container flink-jobmanager-task4 Created 0.1s
  ✓ Container flink-taskmanager-task4 Created 0.2s
  ✓ Container kafka-task4 Created 0.2s
Attaching to flink-jobmanager-task4, flink-taskmanager-task4, kafka-task4, postgres-task4, zookeeper-tas
k4
zookeeper-task4 | ==> User
zookeeper-task4 | uid=1000(appuser) gid=1000(appuser) groups=1000(appuser)
zookeeper-task4 | ==> Configuring ...
flink-jobmanager-task4 | sed: can't read /config.yaml: No such file or directory
kafka-task4 | ==> User
postgres-task4 | The files belonging to this database system will be owned by user "postgres".
```

Step 2 — Data and Kafka Topics

Create Topics

```
Task_4 — -zsh — 141x33
~/MSC/BigData/Assignment/Final2/Task_4 — docker-compose - docker compose up ... -/Documents/MS/BigData/Assignment/Final2/Task_4 — -zsh

lolitha@Lolithas-MacBook-Air Task_4 % docker exec -it kafka-task4 \
kafka-topics \
  --bootstrap-server kafka-task4:9092 \
  --list

lolitha@Lolithas-MacBook-Air Task_4 % docker exec -it kafka-task4 \
kafka-topics \
  --bootstrap-server kafka-task4:9092 \
  --create \
  --topic twitter.topic \
  --partitions 1 \
  --replication-factor 1
[
  WARNING: Due to limitations in metric names, topics with a period ('.') or underscore ('_') could collide. To avoid issues it is best to use
  either, but not both.
  Created topic twitter.topic.
]
lolitha@Lolithas-MacBook-Air Task_4 % docker exec -it kafka-task4 \
kafka-topics \
  --bootstrap-server kafka-task4:9092 \
  --create \
  --topic tiktok.topic \
  --partitions 1 \
  --replication-factor 1
[
  WARNING: Due to limitations in metric names, topics with a period ('.') or underscore ('_') could collide. To avoid issues it is best to use
  either, but not both.
  Created topic tiktok.topic.
]
lolitha@Lolithas-MacBook-Air Task_4 % docker exec -it kafka-task4 \
kafka-topics \
  --bootstrap-server kafka-task4:9092 \
  --list
[
  tiktok.topic
  twitter.topic
]
lolitha@Lolithas-MacBook-Air Task_4 %
```

Send Data To Topic

```
datasets — -zsh — 111x11
...docker-compose - docker compose up ... ...opic twitter.topic --from-beginning ... ..Final2/Task_4/datasets — -zsh

Last login: Thu Jan 1 21:23:52 on ttys008
[ls]
lolitha@Lolithas-MacBook-Air Task_4 % ls
datasets          docker-compose.yml  setup.txt
lolitha@Lolithas-MacBook-Air Task_4 % cd datasets
lolitha@Lolithas-MacBook-Air datasets % tail -n +2 "Twitter- datasets.csv" | docker exec -i kafka-task4 \
kafka-console-producer \
  --bootstrap-server kafka-task4:9092 \
  --topic twitter.topic
[
lolitha@Lolithas-MacBook-Air datasets %
```

Read Data From Topic

```
Task_4 -- docker exec -it kafka-task4 kafka-console-consumer --bootstrap-server kafka-task4:9092 --topic twitter.topic --from-beginning...

Last login: Thu Jan 1 21:14:44 on ttys000
lolitha@lolithas-MacBook-Air Task_4 % docker exec -it kafka-task4 \
  kafka-console-consumer \
    --bootstrap-server kafka-task4:9092 \
    --topic twitter.topic \
    --from-beginning
*1868428607451799983", "GloboNews", "GloboNews", "Com o fim da ditadura Assad, muitos sírios consideram voltar para casa. 12 milhões fugir
am desde o início da guerra civil, há 13 anos. Países europeus estão revendo políticas de refúgio a sírios.

- Assista à #GloboNews: ""2024-12-15T22:03:08.000Z""", "null", "https://x.com/GloboNews/status/1868428607451799983", "null", 2.1.33.8369, "http:
//glo.bo/39WjXAN", [{"GloboNews": {"GloboNewsInternacional": 3771758, "Nunca desliga": 222223, "https://pbs.twimg.com/profile_images/155918271
243933185/evnVnns_normal.jpg", 122, false, 1.1, [{"post_id": "null", "profile_id": "null", "profile_name": "null", "url": "https://video.twimg.com/ext_tw
video_url": "https://video.twimg.com/amplify_video/1258428486872754176/vid/avc1/1928a18380/mba1ZCAdyOxwfoe.mp4?tag=16"}], [{"data_posted":
null, "description": "null", "photos": "null", "post_id": "null", "profile_id": "null", "profile_name": "null", "url": "null", "videos": "null"}]
*1868159894567121215", "billboard", "billboard", "Brian Austin Green Tells MxK to 'Grow Up' After Musician's Split From Megan Fox: 'She's
Pregnant'", ""2024-12-15T08:00:11.000Z""", "null", "https://x.com/billboard/status/1868159894567121215", "null", 7.3.43.21987, "https://www.billbo
ard.com/music/music-news/brian-austin-green-reacts-megan-fox-machine-gun-kelly-split-1235856363/?utm_campaign=trusant&utm_medium=social&ut
m_source=twitter", "null", 14257949, @billboardcharts, @billboardhiphop, @billboardlatin, @billboardpro, @hopbillboard,
@billboarddance, @billboardpride, 367584, "https://pbs.twimg.com/profile_images/16965772024966667/FhNkxwF_normal.jpg", 3784, false, 1.2, [{"
post_id": "null", "profile_id": "null", "profile_name": "null", "url": "https://pbs.twimg.com/card_img/1867636129663848992/-o2NDwaA7formstrjpp&name=s
rig"}], [{"data_posted": "null", "description": "null", "photos": "null", "post_id": "null", "profile_id": "null", "profile_name": "null", "url":
null, "videos": "null"}]
*1868451534708883739", "TNTSports", "TNTSports", "VENCE O PSG NO CLÁSSICO! 🏆 Nossa @clsalbuquerque traz os detalhes de mais um triunfo
de time da capital, dessa vez contra o Lyon, pela #Ligue1 e destaque a atuação de uma jovem promessa da equipe comandada por Luis Enrique!",
""2024-12-16T00:22:14.000Z""", "null", "https://x.com/TNTSportsBR/status/1868451534708883739", [{"biography": "null", "followers": "null", "follow
ing": "null", "is_verified": "null", "profile_id": "29627623", "profile_name": "Clara Albuquerque", "url": "https://x.com/clsalbuquerque"}]
*2.1.33.15497", [{"Ligue1": "#293189", "🏆 é muito mais Champions na Mx, com TODOS os jogos ao vivo!
Assine e assista clicando no link 📡, 456734, "https://pbs.twimg.com/profile_images/1887013047975768865/8sm_i0pw_normal.jpg", 859, false, 8.
1, [{"post_id": "null", "profile_id": "null", "profile_name": "null", "url": "https://video.twimg.com/ext_tw
video/1868451336026205549/pu/vid/avc1/720x1280/tw3z5XZmhVY4hfX.mp4?tag=12"}], [{"data_posted": "null", "description": "null", "photos": "null",
"post_id": "null", "profile_id": "null", "profile_name": "null", "url": "null", "videos": "null"}]
*186844138202717466", "TNTSports", "TNTSports", "IDOLLO É AGORA PRESIDENTE! 🇲🇵 ex-atacante Diego Milito foi confirmado como novo pres
idente do Racing!

@TNTSportsAR", ""2024-12-15T23:41:04.000Z""", [{"https://pbs.twimg.com/media/Ge43hOgKcAAQy5K.jpg"}], "https://x.com/TNTSportsBR/status/1868
44138202717466", [{"biography": "null", "followers": "null", "following": "null", "is_verified": "null", "profile_id": "854784383982723872", "pro
file_name": "TNT Sports Argentina", "url": "https://x.com/TNTSportsAR"}], 6.5.238.18267, "null", 8293189, "🏆 é muito mais Champions na M
x, com TODOS os jogos ao vivo!
Assine e assista clicando no link 📡, 456734, "https://pbs.twimg.com/profile_images/1887013047975768865/8sm_i0pw_normal.jpg", 859, false, 1.
3, [{"post_id": "null", "profile_id": "null", "profile_name": "null", "url": "https://video.twimg.com/ext_tw
video/1868451336026205549/pu/vid/avc1/720x1280/tw3z5XZmhVY4hfX.mp4?tag=12"}], [{"data_posted": "null", "description": "null", "photos": "null",
"post_id": "null", "profile_id": "null", "profile_name": "null", "url": "null", "videos": "null"}]
*1868418260892565925", "GloboNews", "GloboNews", ".@DanielLima : cirurgião de Lula travou negociações para reforma ministerial, que ficou
```

Step 3 — Flink Stream Processing

Flink Streaming Application Implementation

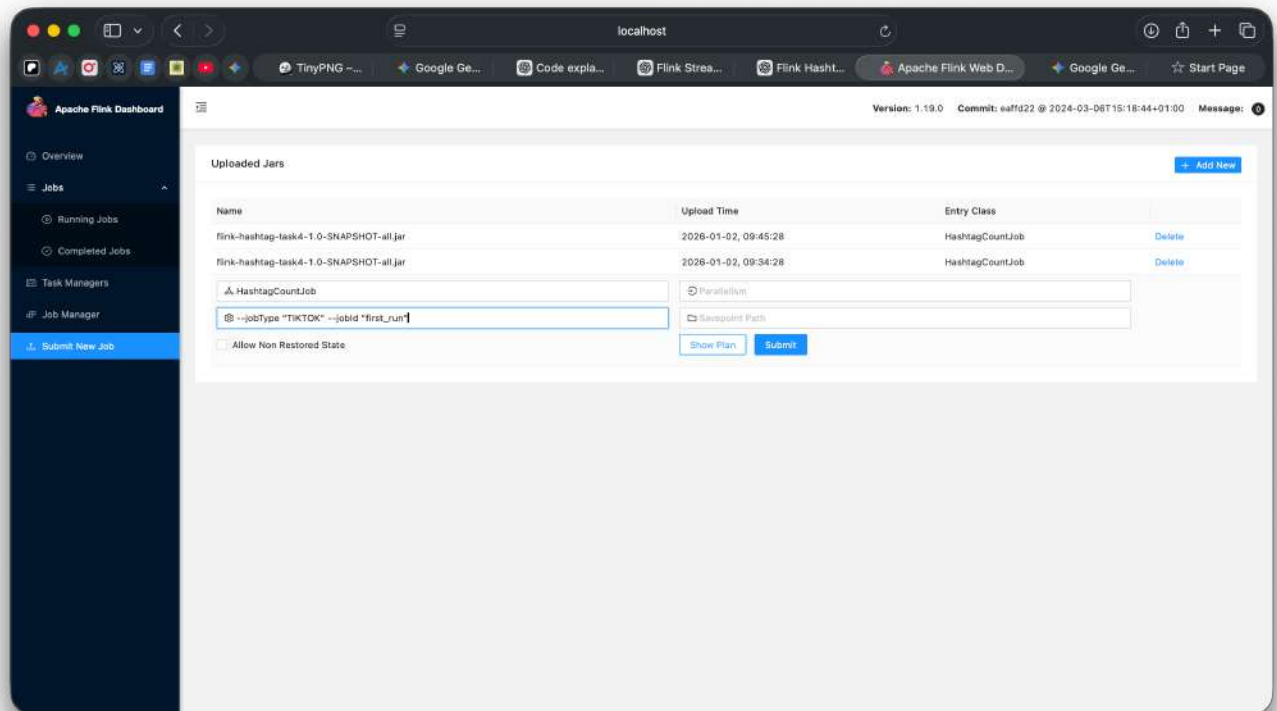
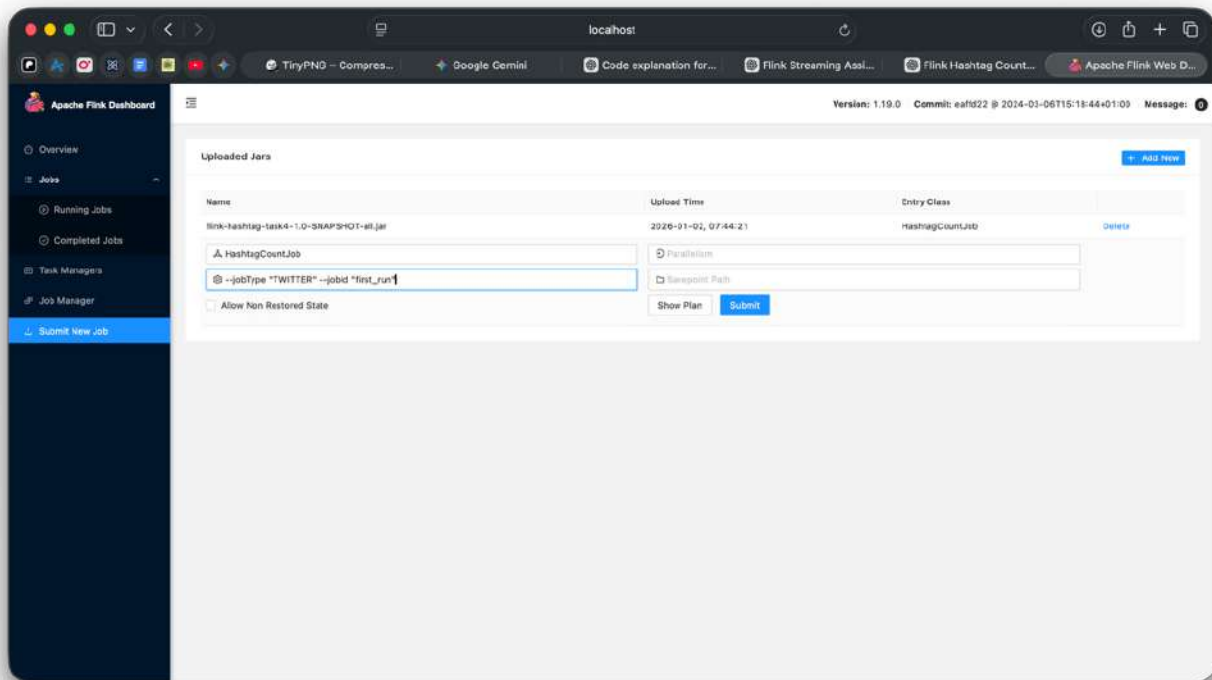
A Flink Stream Application was developed using Java to handle Twitter and TikTok datasets which would count the hashtags #GloboNews and #SAE.

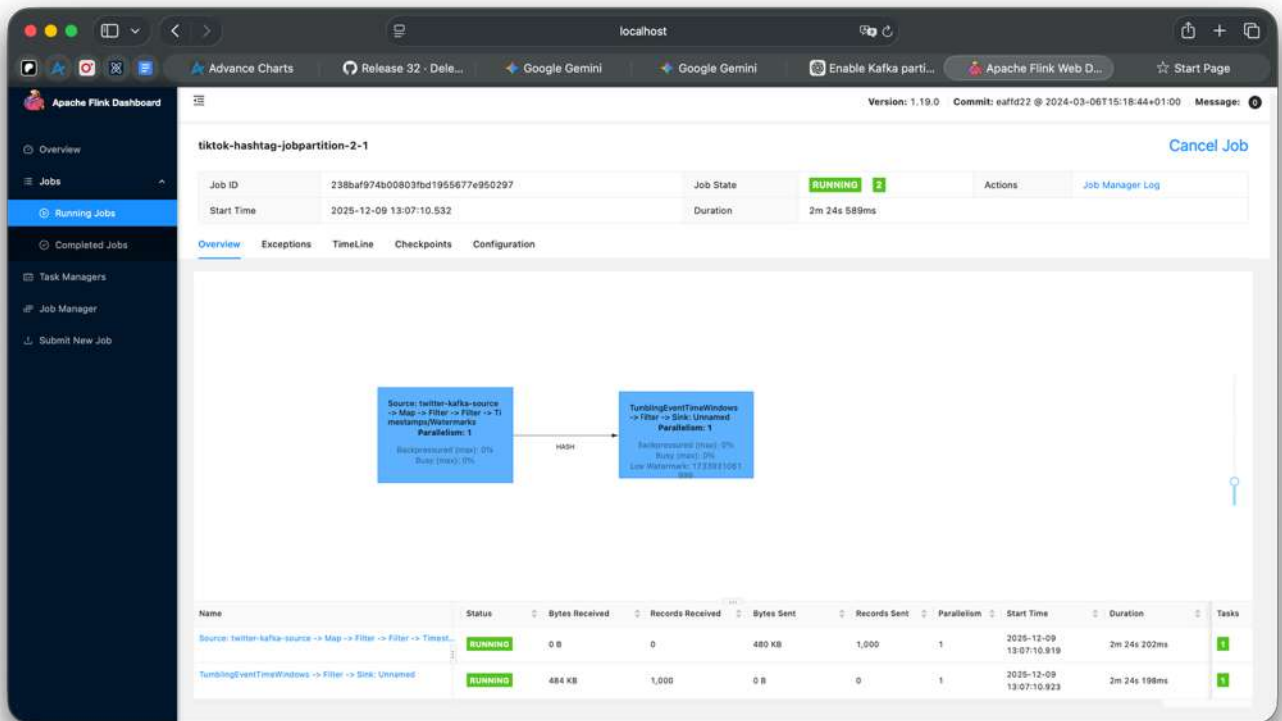
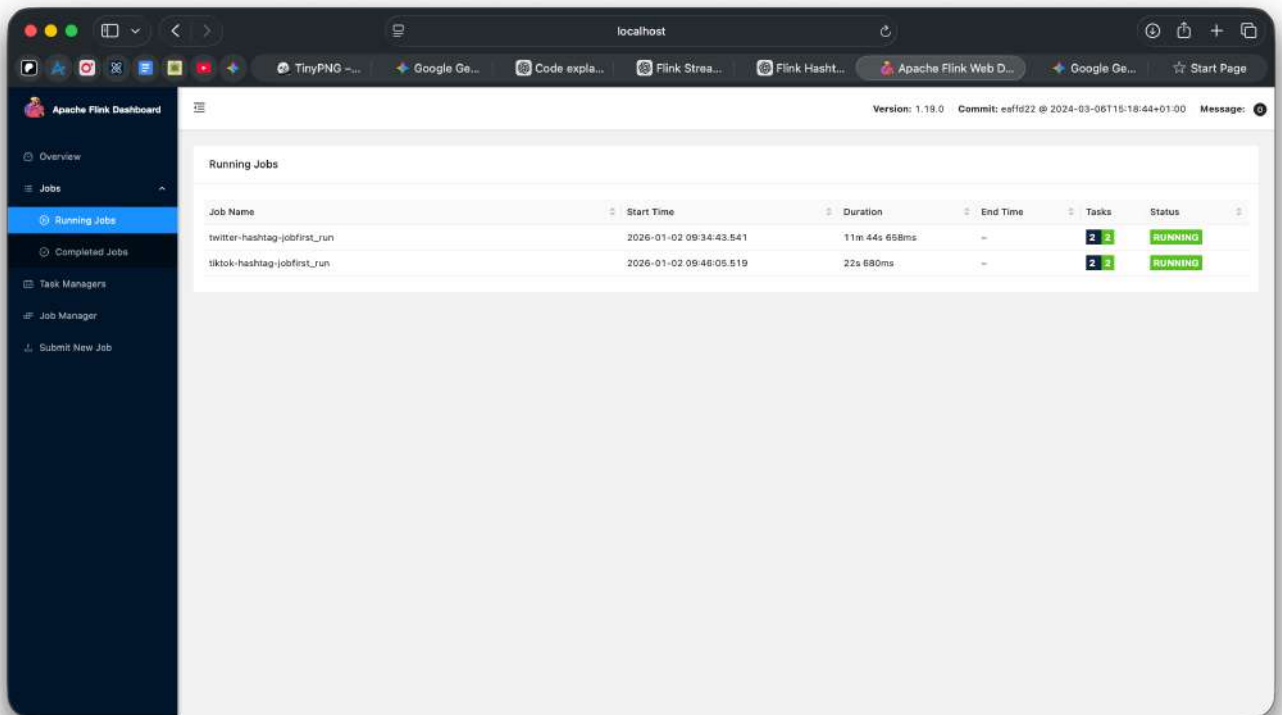
Two Kafka topics were created, one carrying Twitter CSV records and the other carrying TikTok CSV records. “kafka-console-producer” was used to push the records to the kafka topics. The Flink Application consume the messages using the “flink-connector-kafka” dependency and then process the data to output results to the PostgreSQL.

Heart of the Flink Application is “HashtagCountJob” which carries out the following duties:

- Creates the StreamExecutionEnvironment.
- Reads command line args to decide which dataset to run on (Twitter or TikTok) and picks the correct Kafka topic, group id, hashtag, and job name.
- Then it reads Kafka messages as plain strings using the Kafka source (SimpleStringSchema) from "flink-connector-kafka" lib.
- EventParser class identifies and extract the event time using a regex and parses each string into an Event.
- It assigns event time and then HashtagCountJob assign watermarks.
- HashtagCountJob groups the events and runs a tumbling event time window.
- The window length is controlled by WINDOW_SECONDS = 15. The job uses bounded out of order watermarks with a tolerance of WATERMARK_DELAY_SECONDS = 5. Therefore, Flink allows events to arrive up to 5 seconds late before a window is considered complete. The job does not configure allowedLateness. So, events arriving after the watermark passes the window end will not be included in that window's counts.
- For each window it outputs a Result that contains rawEventCount and hashtagEventCount.
- Finally, it saves the results to PostgreSQL using via a JDBC sink.
- The performance is captured using a small script which will record the memory and cpu usage every 1 second using docker stats command for the “flink-taskmanager-task4” container.
- Source : https://github.com/lolitha-lakshan-1/bigdata-assignment/tree/main/Task_4/flink-hashtag-count-gradle/src/main/java

Flink Job Submission





Result Accuracy and Outputs

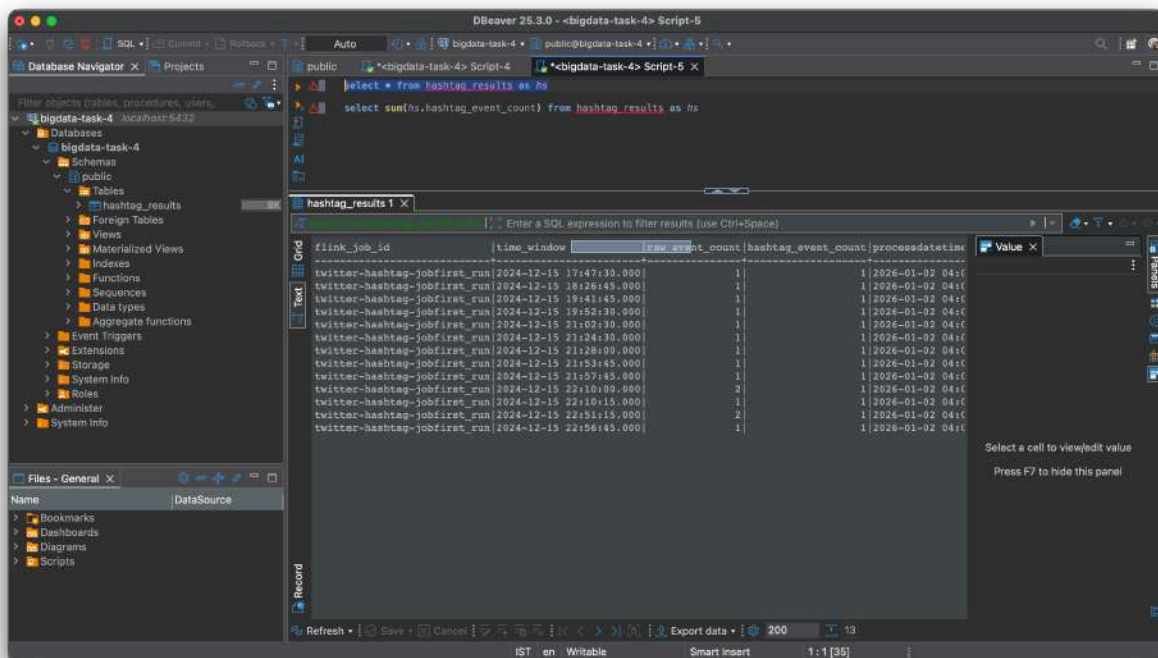
Twitter Dataset

Expected “#GloboNews” Hashtag Count: 13



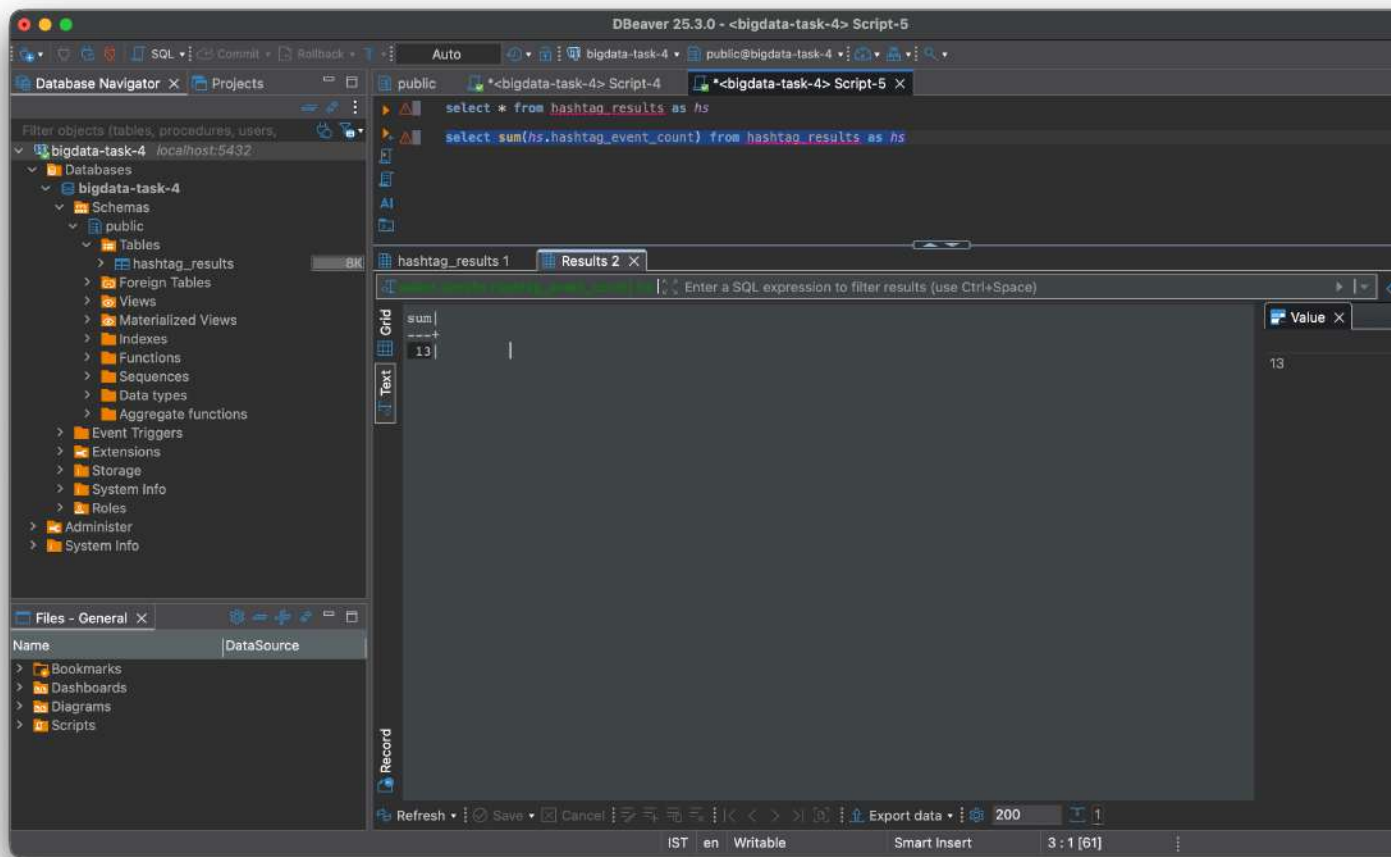
```
loliitha@Loliithas-MacBook-Air datasets % grep -c "#GloboNews" Twitter-datasets.csv
13
loliitha@Loliithas-MacBook-Air datasets %
```

Flink Output :



The screenshot shows the DBBeaver interface with a Flink job named 'hashtag_results' running. The job is configured with a 'select * from hashtag_results as hs' query. The output table, 'hashtag_results', displays the following data:

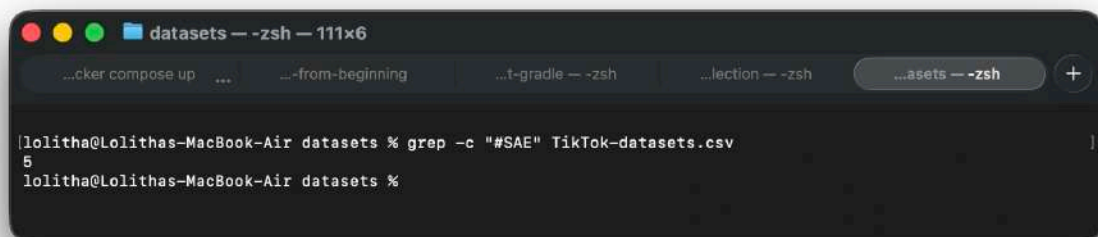
| flink_job_id | time_window | raw_event_count | hashtag_event_count | processdatetime |
|------------------------------|-------------------------|-----------------|---------------------|---------------------|
| twitter-hashtag-jobfirst_run | 2024-12-15 17:47:30.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 18:26:45.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 19:41:45.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 19:52:30.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 21:02:30.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 21:24:30.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 21:28:00.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 21:53:45.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 21:57:45.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 22:10:00.000 | 2 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 22:10:15.000 | 1 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 22:51:15.000 | 2 | 1 | 2026-01-02 04:00:00 |
| twitter-hashtag-jobfirst_run | 2024-12-15 22:56:45.000 | 1 | 1 | 2026-01-02 04:00:00 |



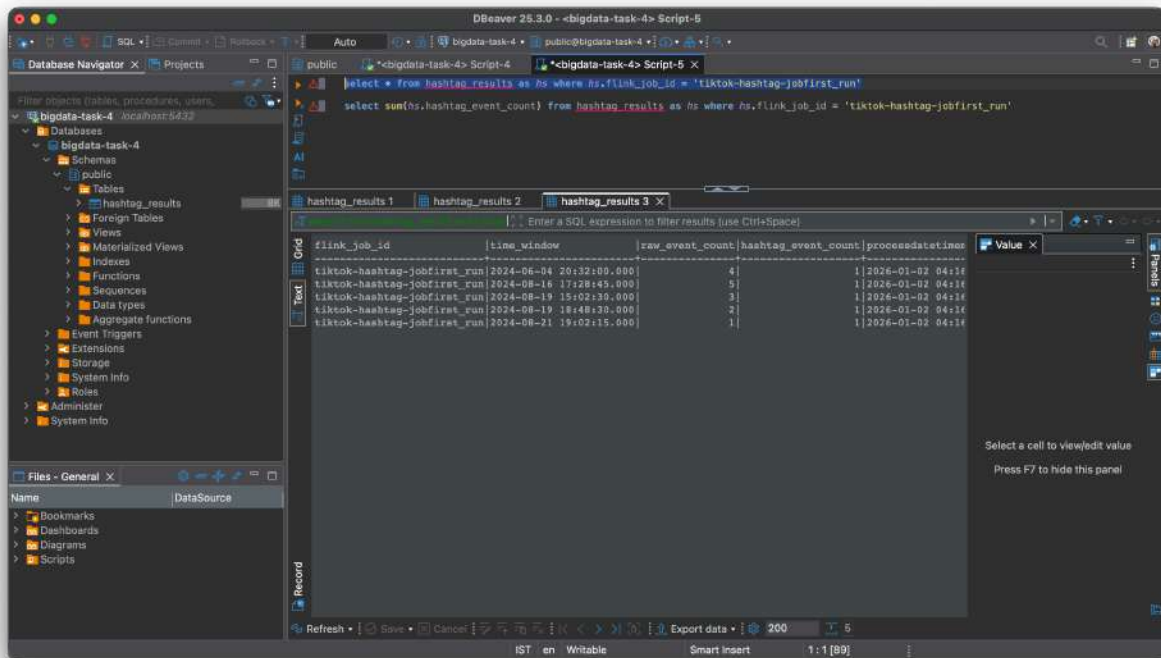
“#GloboNews” hashtag count is 13 therefore, result is accurate.

TikTok Dataset

Expected “#SAE” Hashtag Count: 5



Flink Output:

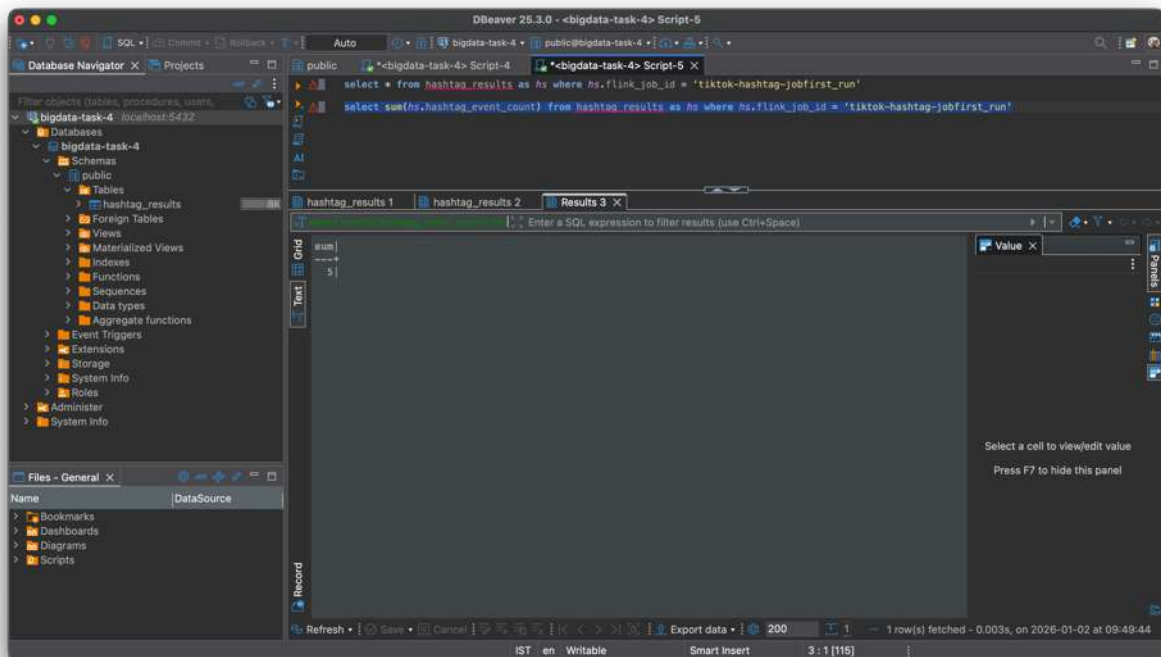


The screenshot shows the DBeaver 25.3.0 interface. The SQL editor contains the following query:

```
select * from hashtag_results as hs where hs.flink_job_id = 'tiktok-hashtag-jobfirst_run'
```

The results are displayed in a grid view with the following columns: flink_job_id, time_window, raw_event_count, hashtag_event_count, and processed_at_time. The data is as follows:

| flink_job_id | time_window | raw_event_count | hashtag_event_count | processed_at_time |
|-----------------------------|-------------------------|-----------------|---------------------|-------------------|
| tiktok-hashtag-jobfirst_run | 2024-06-04 20:32:00.000 | 4 | 1 | 2026-01-02 04:16 |
| tiktok-hashtag-jobfirst_run | 2024-06-16 17:28:45.000 | 5 | 1 | 2026-01-02 04:16 |
| tiktok-hashtag-jobfirst_run | 2024-06-19 15:02:30.000 | 3 | 1 | 2026-01-02 04:16 |
| tiktok-hashtag-jobfirst_run | 2024-06-19 16:48:30.000 | 2 | 1 | 2026-01-02 04:16 |
| tiktok-hashtag-jobfirst_run | 2024-06-21 15:02:15.000 | 2 | 1 | 2026-01-02 04:16 |



The screenshot shows the same DBeaver 25.3.0 interface, but the results are now displayed in a summary view. The query is:

```
select * from hashtag_results as hs where hs.flink_job_id = 'tiktok-hashtag-jobfirst_run'
```

```
select sum(hs.hashtag_event_count) from hashtag_results as hs where hs.flink_job_id = 'tiktok-hashtag-jobfirst_run'
```

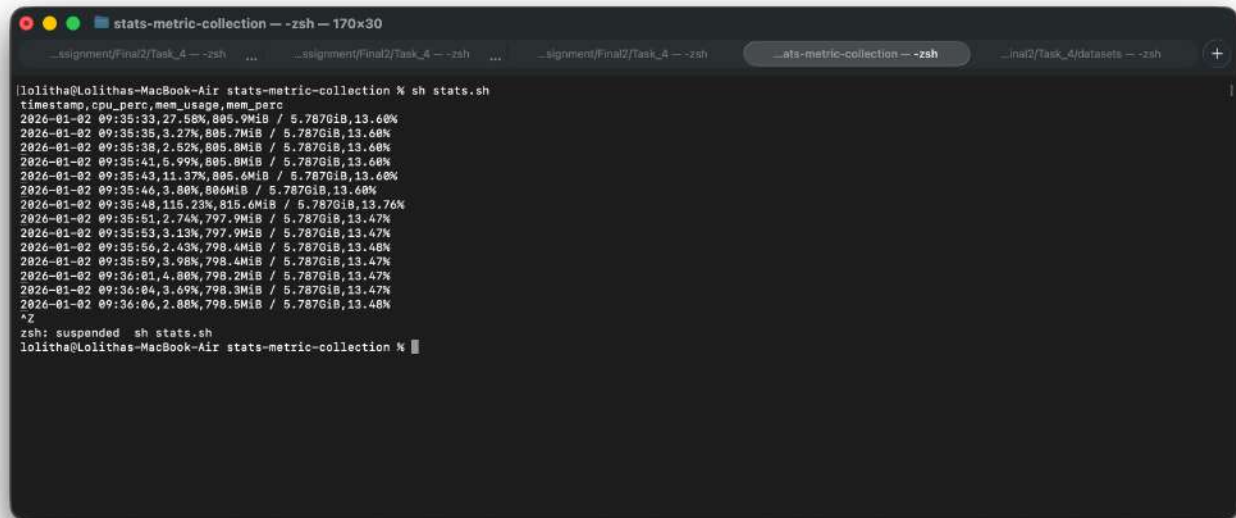
The summary view shows a single row with the value 5 for the sum of hashtag_event_count.

| sum |
|-----|
| 5 |

“#SAE” hashtag count is 5 therefore, result is accurate.

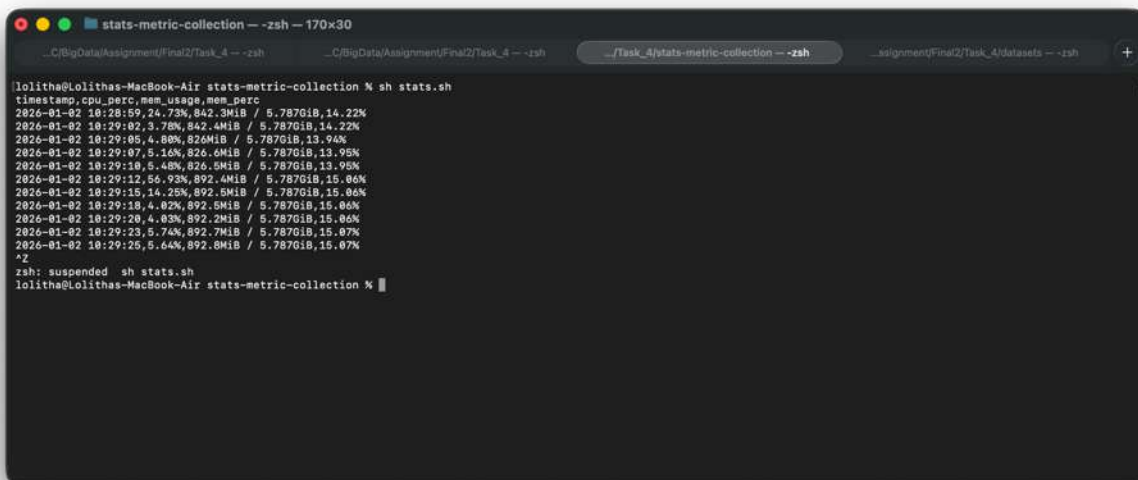
Performance

Metrics for Twitter Job

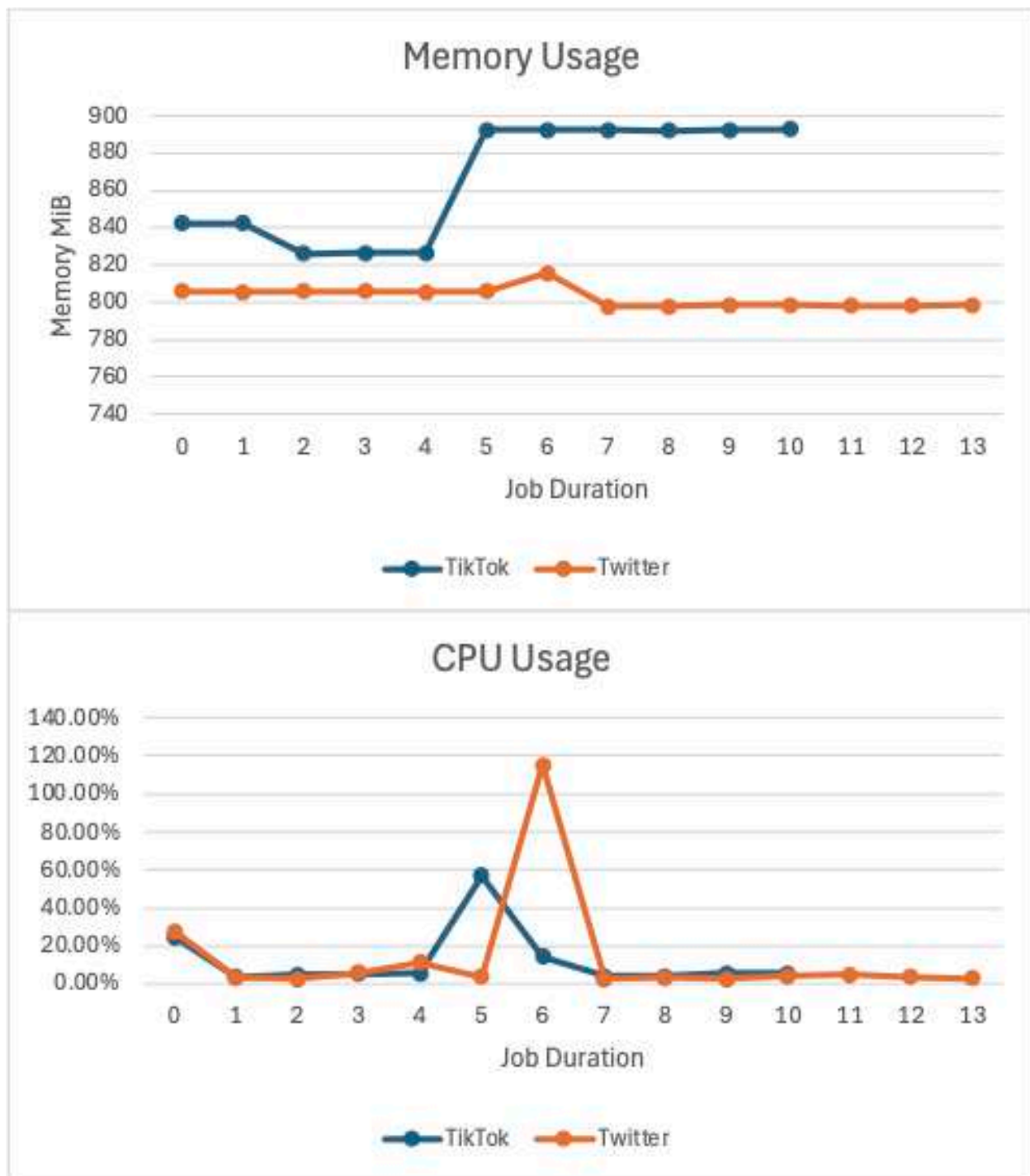
A terminal window titled 'stats-metric-collection --zsh -- 170x30' showing the output of a script 'stats.sh'. The script prints a table of system metrics (timestamp, cpu_perc, mem_usage, mem_perc) for a duration of 10 seconds. The metrics are collected every 1 second. The output shows that the system is running on a MacBook-Air with 5.787GiB of memory and 13.60% of memory usage. The CPU usage is around 3.27% to 3.69%. The script is suspended and then resumed.

```
lolitha@Lolithas-MacBook-Air stats-metric-collection % sh stats.sh
timestamp,cpu_perc,mem_usage,mem_perc
2026-01-02 09:35:33,27.58%,805.9MiB / 5.787GiB,13.60%
2026-01-02 09:35:34,3.27%,805.7MiB / 5.787GiB,13.60%
2026-01-02 09:35:35,3.27%,805.7MiB / 5.787GiB,13.60%
2026-01-02 09:35:36,2.52%,805.8MiB / 5.787GiB,13.60%
2026-01-02 09:35:37,5.99%,805.8MiB / 5.787GiB,13.60%
2026-01-02 09:35:38,11.37%,805.6MiB / 5.787GiB,13.60%
2026-01-02 09:35:39,3.80%,804MiB / 5.787GiB,13.60%
2026-01-02 09:35:40,115.23%,815.6MiB / 5.787GiB,13.76%
2026-01-02 09:35:41,2.74%,797.9MiB / 5.787GiB,13.47%
2026-01-02 09:35:42,3.13%,797.9MiB / 5.787GiB,13.47%
2026-01-02 09:35:43,2.43%,798.4MiB / 5.787GiB,13.48%
2026-01-02 09:35:44,3.98%,798.4MiB / 5.787GiB,13.47%
2026-01-02 09:35:45,4.80%,798.2MiB / 5.787GiB,13.47%
2026-01-02 09:35:46,3.69%,798.3MiB / 5.787GiB,13.47%
2026-01-02 09:35:47,2.88%,798.5MiB / 5.787GiB,13.48%
^Z
zsh: suspended sh stats.sh
lolitha@Lolithas-MacBook-Air stats-metric-collection %
```

Metrics for TikTok Job

A terminal window titled 'stats-metric-collection --zsh -- 170x30' showing the output of a script 'stats.sh'. The script prints a table of system metrics (timestamp, cpu_perc, mem_usage, mem_perc) for a duration of 10 seconds. The metrics are collected every 1 second. The output shows that the system is running on a MacBook-Air with 5.787GiB of memory and 14.22% of memory usage. The CPU usage is around 3.78% to 5.64%. The script is suspended and then resumed.

```
lolitha@Lolithas-MacBook-Air stats-metric-collection % sh stats.sh
timestamp,cpu_perc,mem_usage,mem_perc
2026-01-02 10:28:59,24.73%,842.3MiB / 5.787GiB,14.22%
2026-01-02 10:29:00,3.78%,842.4MiB / 5.787GiB,14.22%
2026-01-02 10:29:01,4.80%,826MiB / 5.787GiB,13.94%
2026-01-02 10:29:02,5.16%,826.6MiB / 5.787GiB,13.95%
2026-01-02 10:29:03,5.48%,826.5MiB / 5.787GiB,13.95%
2026-01-02 10:29:04,56.93%,892.4MiB / 5.787GiB,15.06%
2026-01-02 10:29:05,14.25%,892.5MiB / 5.787GiB,15.06%
2026-01-02 10:29:06,4.02%,892.5MiB / 5.787GiB,15.06%
2026-01-02 10:29:07,4.03%,892.2MiB / 5.787GiB,15.06%
2026-01-02 10:29:08,5.74%,892.7MiB / 5.787GiB,15.07%
2026-01-02 10:29:09,5.64%,892.8MiB / 5.787GiB,15.07%
^Z
zsh: suspended sh stats.sh
lolitha@Lolithas-MacBook-Air stats-metric-collection %
```



Twitter dataset size is around 1 MB and TikTok dataset size is around 500 kilobytes. Flink jobs completed in a span of 2 – 3 seconds after receiving data. Containers were restarted before each job, and it was noted that Flink allocated more base memory when it came to the tiktok job. The Tiktok job recorded a ~+60 MiB spike of memory compared to ~+20 MiB spike of Twitter Job. Furthermore, on both occasions there were CPU spikes observed, Twitter Job recorded a spike 2X compared to TikTok.

Step 4 — Scaling Experiment

Implementation

The parallel Flink job is implemented in “HashtagCountParallelJob” class which carries out the following duties:

- It creates the `StreamExecutionEnvironment` and reads command line args to decide which dataset to run on (Twitter or TikTok) and picks the correct Kafka topic, group id, hashtag, and job name.
- It uses the two new partition topics to support higher parallel consumption. The topics are “twitter.topic-partition-2” and “tiktok.topic-partition-2” created with partitions size of 2.
- Increased Flink operator parallelism by calling `environment.setParallelism(2)`. This allows the Kafka source to consume the two topic partitions in parallel.
- It reads Kafka messages as plain strings using the Kafka source with `SimpleStringSchema` from the `flink-connector-kafka` library and assigned watermarks at the Kafka source. On the single partition job this was done after passing the event.
- The watermark strategy is bounded out of order and uses watermark delay of 5 seconds. The timestamp assigner extracts event time by parsing each record with `EventParser.parse(record)` and returning `event.getEventTimeMillis()` and filters any invalid timestamps.
- It groups events and runs a tumbling event time window. The window length is controlled by `WINDOW_SECONDS = 15`. The job does not configure `allowedLateness`. So events arriving after the watermark passes the window end will not be included in that window’s counts.
- For each window it outputs a `Result` that contains `rawEventCount` and `hashtagEventCount`.
- Finally, it saves the results to PostgreSQL using a JDBC sink.
- The performance is captured using the same docker stats script which records CPU and memory usage every 1 second for the Flink containers.

Source : https://github.com/lolitha-lakshan-1/bigdata-assignment/blob/main/Task_4/flink-hashtag-count-gradle/src/main/java/HashtagCountParallelJob.java

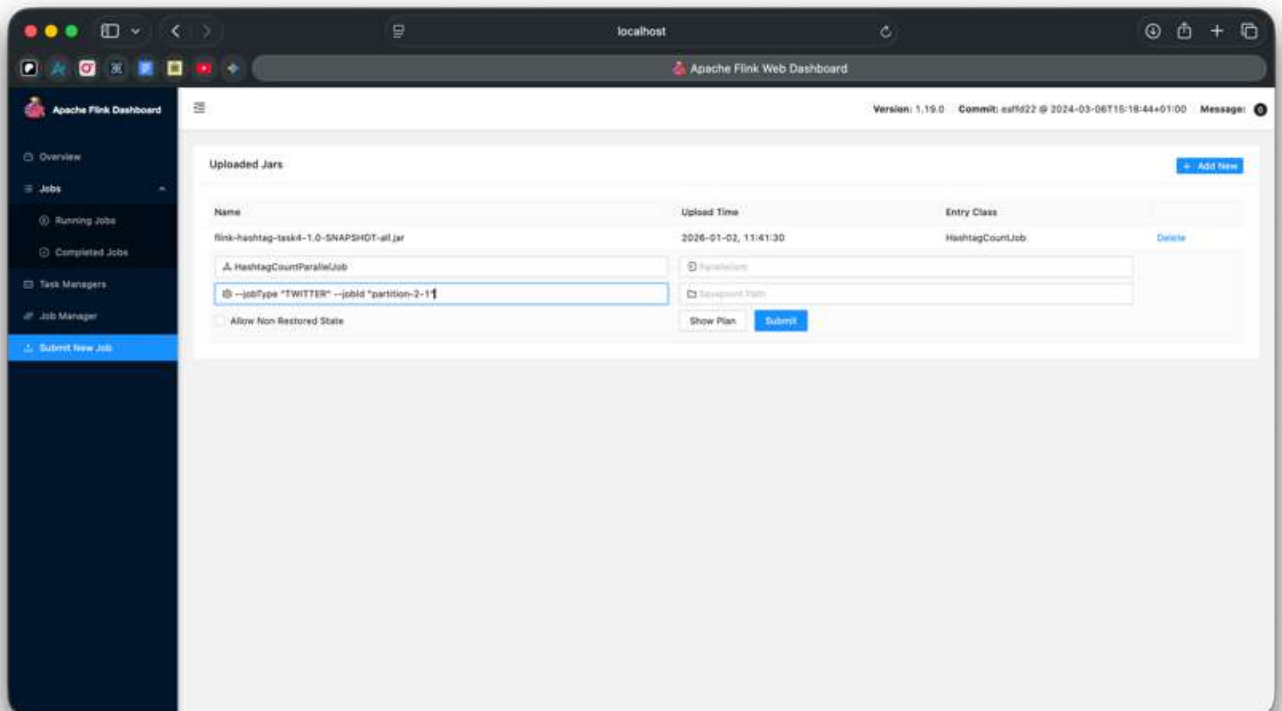
```
stats-metric-collection --zsh -- 170x22
...ignment/Final2/Task_4 -- docker-compose - docker compose up ... ...ata/Assignment/Final2/Task_4/stats-metric-collection -- -zsh ...signment/Final2/Task_4/stats-metric-collection -- -zsh

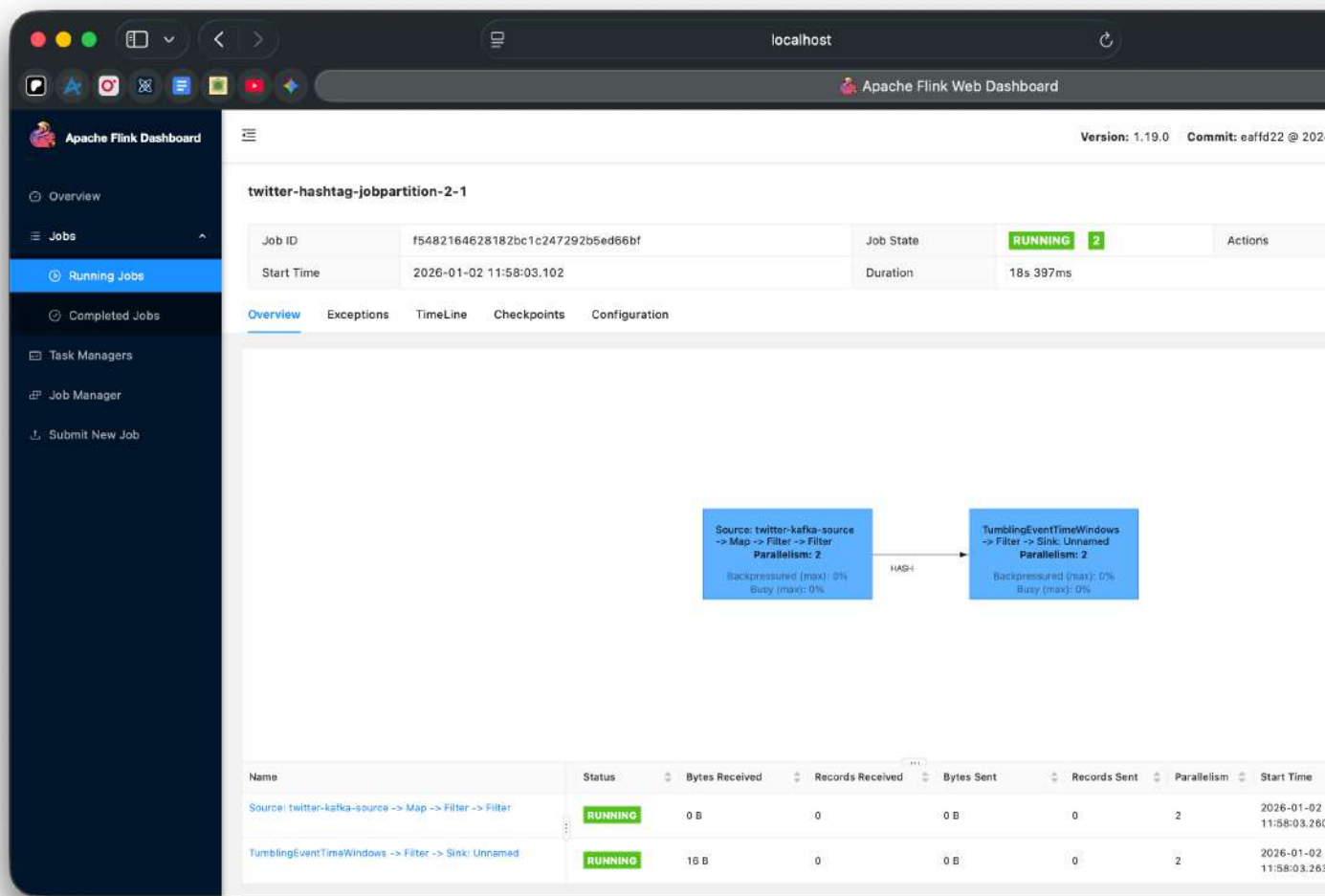
Last login: Fri Jan  2 10:23:27 on ttys000
lolitha@Lolithas-MacBook-Air stats-metric-collection % docker exec -it kafka-task4 \
kafka-topics \
--bootstrap-server kafka-task4:9092 \
--list

lolitha@Lolithas-MacBook-Air stats-metric-collection % docker exec -it kafka-task4 \
kafka-topics \
--bootstrap-server kafka-task4:9092 \
--create \
--topic twitter.topic-partition-2 \
--partitions 2 \
--replication-factor 1

WARNING: Due to limitations in metric names, topics with a period ('.') or underscore ('_') could collide. To avoid issues it is best to use either, but not both.
Created topic twitter.topic-partition-2.
lolitha@Lolithas-MacBook-Air stats-metric-collection % docker exec -it kafka-task4 \
kafka-topics \
--bootstrap-server kafka-task4:9092 \
--list
twitter.topic-partition-2
lolitha@Lolithas-MacBook-Air stats-metric-collection %
```

Flink Job Submission





Now Flink Job indicates the parallelism of 2.

Result Accuracy and Outputs

Twitter Dataset

Expected “#GloboNews” Hashtag Count: 13

Even with parallelism enabled Flink returned hashtag count of 13 for “#GloboNews” hashtag.

Therefore, result is accurate.

Flink Output:

DBeaver 25.3.0 - <bigdata-task-4> Script-4

```

-- bigdata-task-4 Script-4
--
-- Schema: public
--
-- Table: hashtag_results
--
-- Columns:
--   flink_job_id TEXT NOT NULL,
--   time_window TIMESTAMP NOT NULL,
--   raw_event_count BIGINT NOT NULL,
--   hashtag_event_count BIGINT NOT NULL,
--   processdatetimeinmillis TIMESTAMP NOT NULL
--
-- SQL:
-- select * from hashtag_results as hs
-- select sum(hs.raw_event_count) from hashtag_results as hs

```

| flink_job_id | time_window | raw_event_count | hashtag_event_count | processdatetimeinmillis |
|----------------------------------|-------------------------|-----------------|---------------------|---------------------------|
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 17:47:30.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 18:26:45.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 19:41:45.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 19:52:30.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 21:02:30.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 21:24:30.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 21:28:00.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 21:53:45.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 21:57:45.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 22:19:00.000 | 1 | 1 | 1/2026-01-02 06:35:44.340 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 22:19:15.000 | 1 | 1 | 1/2026-01-02 06:35:44.341 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 22:51:15.000 | 2 | 1 | 1/2026-01-02 06:35:44.341 |
| twitter-hashtag-jobpartition-2-1 | 2026-12-15 22:56:45.000 | 1 | 1 | 1/2026-01-02 06:35:44.341 |

localhost Apache Flink Web Dashboard

Version: 1.19.0 Commit: eafdd22 @ 2024-03-06T15:18:44+01:00 Message: 0

twitter-hashtag-jobpartition-2-1

Cancel Job

| Job ID | Job State | Actions |
|-----------------------------------|------------------|---------------------------------|
| f5482164628182b0c1c247292b5ed66bf | RUNNING 2 | Job Manager Log |

Start Time: 2026-01-02 11:58:03.102 Duration: 8m 40s 761ms

Overview Exceptions TimeLine Checkpoints Configuration

```

graph LR
    A["Source: twitter-kafka-source  
-> Map -> Filter -> Filter  
Parallelism: 2  
Backpressured (max): N/A  
Busy (max): N/A"] -- HASH --> B["TumblingEventTimeWindows  
-> Filter -> Sink: Unnamed  
Parallelism: 2  
Backpressured (max): N/A  
Busy (max): N/A"]

```

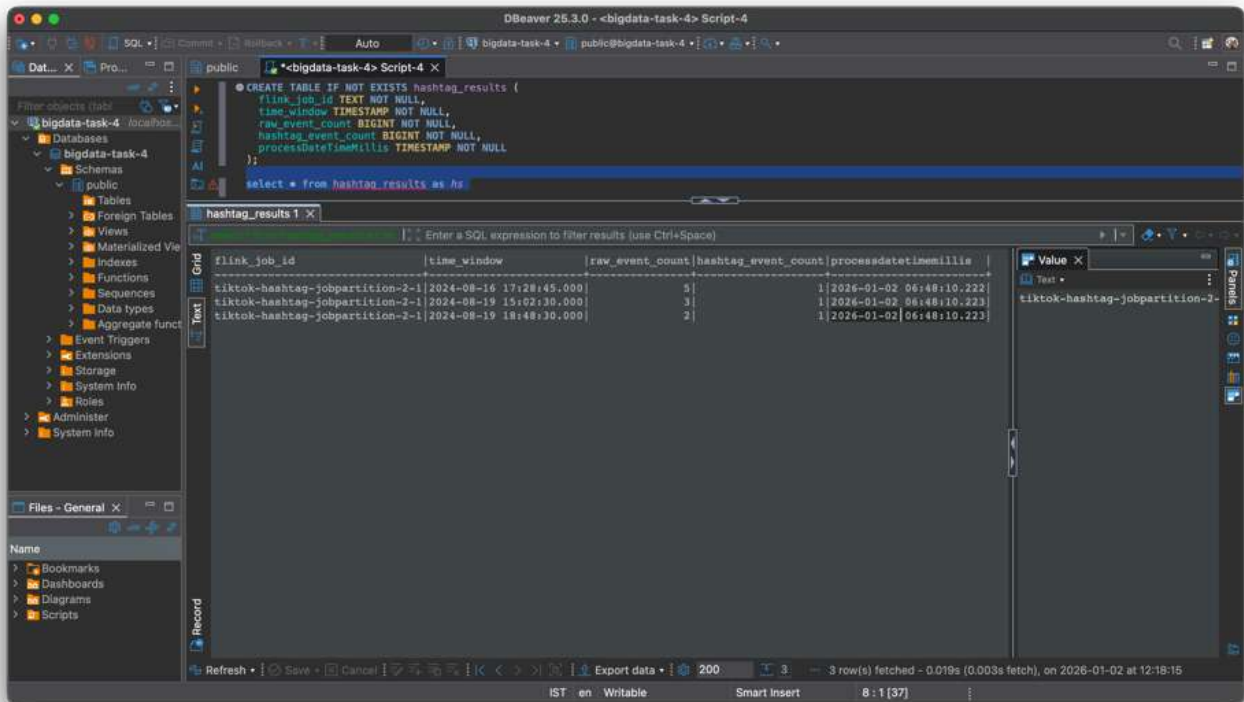
| Name | Status | Bytes Received | Records Received | Bytes Sent | Records Sent | Parallelism | Start Time | Duration | Errors |
|---|----------------|----------------|------------------|------------|--------------|-------------|-------------------------|--------------|--------|
| Source: twitter-kafka-source -> Map -> Filter -> Filter | RUNNING | 0 B | 0 | 672 KB | 1,000 | 2 | 2026-01-02 11:58:03.260 | 8m 40s 603ms | 2 |
| TumblingEventTimeWindows -> Filter -> Sink: Unnamed | RUNNING | 710 KB | 1,000 | 0 B | 0 | 2 | 2026-01-02 11:58:03.263 | 8m 40s 600ms | 2 |

TikTok Dataset

Expected “#SAE” Hashtag Count: 5

However, Flink Parallel job return “#SAE” Hashtag Count: 3 which is incorrect.

Therefore result was inaccurate .



The screenshot shows the DBeaver 25.3.0 interface. The SQL editor contains the following script:

```
CREATE TABLE IF NOT EXISTS hashtag_results (  
    flink_job_id TEXT NOT NULL,  
    time_window TIMESTAMP NOT NULL,  
    raw_event_count BIGINT NOT NULL,  
    hashtag_event_count BIGINT NOT NULL,  
    processDateTimeMillis TIMESTAMP NOT NULL  
);  
  
select * from hashtag_results as hs;
```

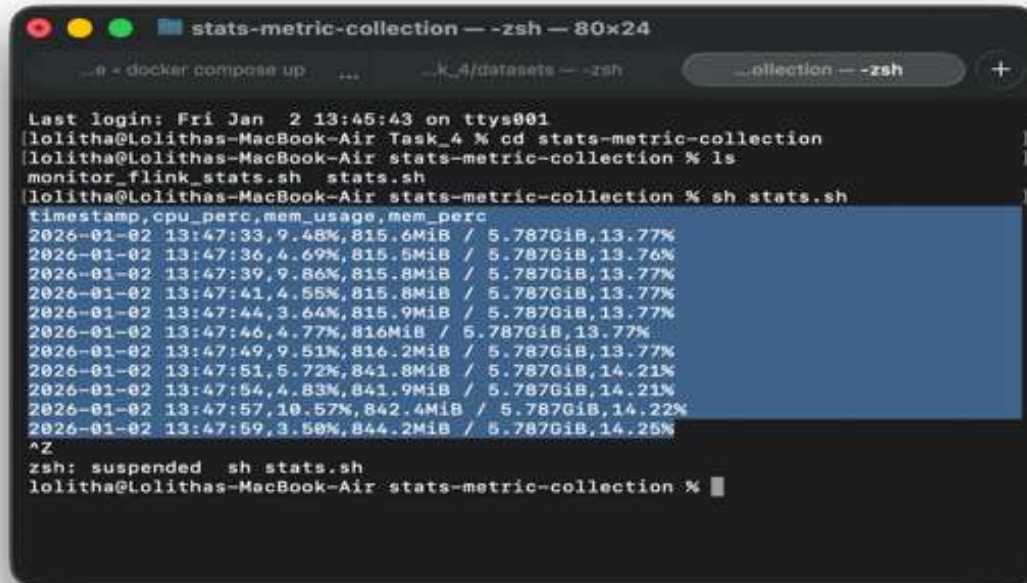
The results pane displays the following data:

| flink_job_id | time_window | raw_event_count | hashtag_event_count | processDateTimeMillis |
|---------------------------------|-------------------------|-----------------|---------------------|-------------------------|
| tiktok-hashtag-jobpartition-2-1 | 2024-08-16 17:28:45.000 | 5 | 1 | 2026-01-02 06:48:10.222 |
| tiktok-hashtag-jobpartition-2-1 | 2024-08-19 15:02:30.000 | 3 | 1 | 2026-01-02 06:48:10.223 |
| tiktok-hashtag-jobpartition-2-1 | 2024-08-19 18:48:30.000 | 2 | 1 | 2026-01-02 06:48:10.223 |

The status bar at the bottom indicates "3 row(s) fetched - 0.019s (0.003s fetch), on 2026-01-02 at 12:18:15".

Performance

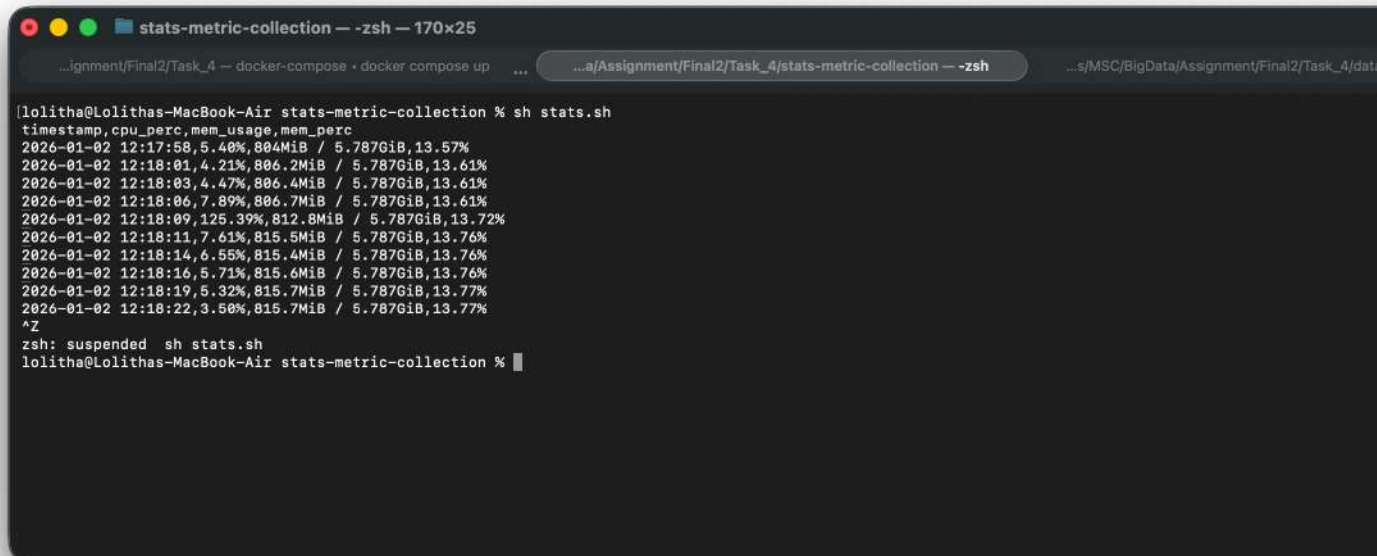
Metrics for Twitter Job



A terminal window titled 'stats-metric-collection --zsh -- 80x24'. The terminal shows the user 'lolitha' at 'Lolithas-MacBook-Air' in the directory 'Task_4'. They run 'cd stats-metric-collection' and 'ls', showing files 'monitor_flink_stats.sh' and 'stats.sh'. Then they run 'sh stats.sh', which outputs a table of metrics. The table has columns: timestamp, cpu_perc, mem_usage, and mem_perc. The data shows metrics for 2026-01-02 from 13:47:33 to 13:47:59. The terminal is then suspended with '^Z' and 'zsh: suspended sh stats.sh'.

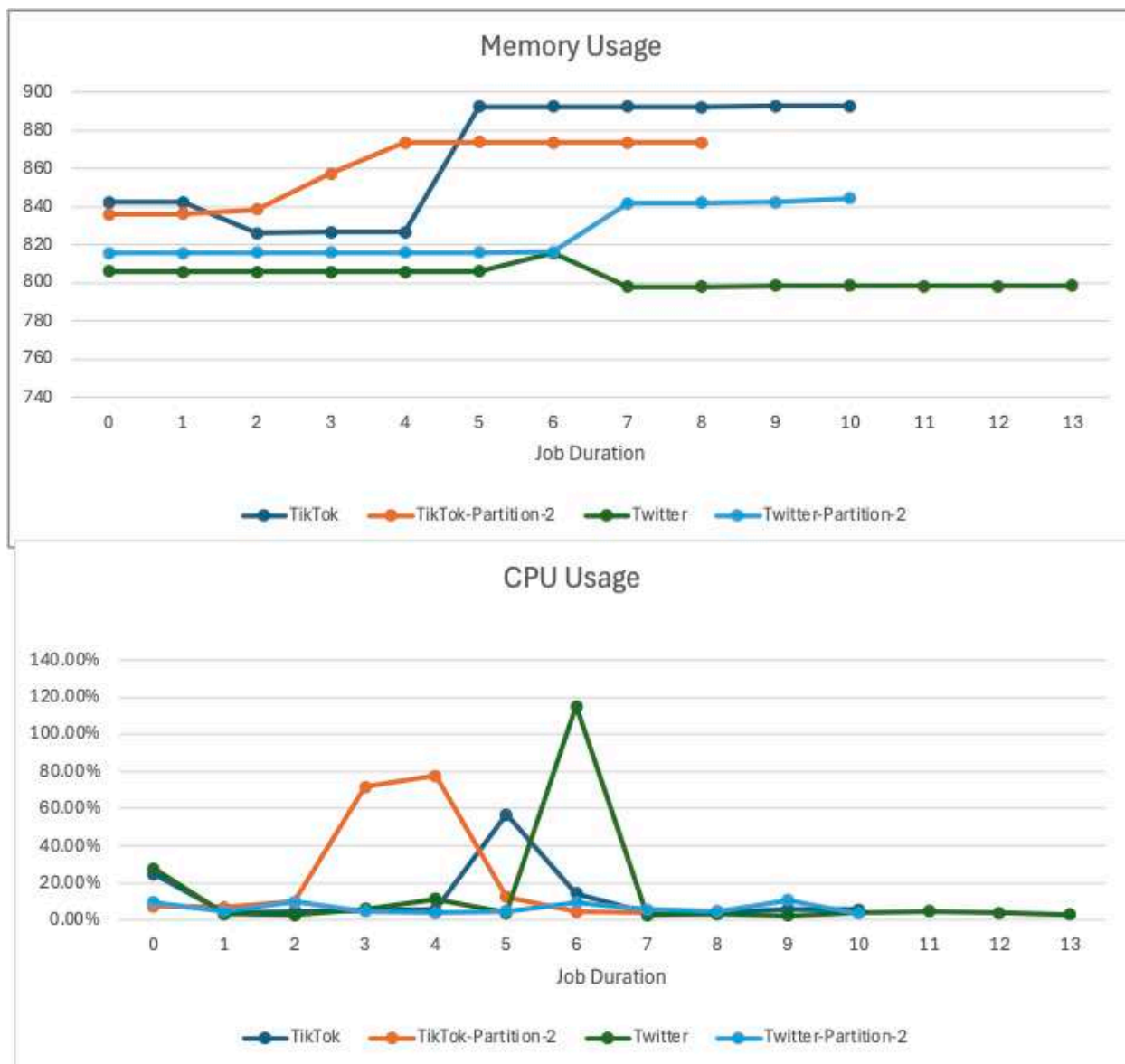
```
Last login: Fri Jan  2 13:45:43 on ttys001
lolitha@Lolithas-MacBook-Air Task_4 % cd stats-metric-collection
lolitha@Lolithas-MacBook-Air stats-metric-collection % ls
monitor_flink_stats.sh  stats.sh
lolitha@Lolithas-MacBook-Air stats-metric-collection % sh stats.sh
timestamp,cpu_perc,mem_usage,mem_perc
2026-01-02 13:47:33,9.48%,815.6MiB / 5.787GiB,13.77%
2026-01-02 13:47:36,4.69%,815.5MiB / 5.787GiB,13.76%
2026-01-02 13:47:39,9.86%,815.8MiB / 5.787GiB,13.77%
2026-01-02 13:47:41,4.55%,815.8MiB / 5.787GiB,13.77%
2026-01-02 13:47:44,3.64%,815.9MiB / 5.787GiB,13.77%
2026-01-02 13:47:46,4.77%,816MiB / 5.787GiB,13.77%
2026-01-02 13:47:49,9.51%,816.2MiB / 5.787GiB,13.77%
2026-01-02 13:47:51,5.72%,841.8MiB / 5.787GiB,14.21%
2026-01-02 13:47:54,4.83%,841.9MiB / 5.787GiB,14.21%
2026-01-02 13:47:57,10.57%,842.4MiB / 5.787GiB,14.22%
2026-01-02 13:47:59,3.50%,844.2MiB / 5.787GiB,14.25%
^Z
zsh: suspended  sh stats.sh
lolitha@Lolithas-MacBook-Air stats-metric-collection %
```

Metrics for TikTok Job



A terminal window titled 'stats-metric-collection --zsh -- 170x25'. The terminal shows the user 'lolitha' at 'Lolithas-MacBook-Air' in the directory 'stats-metric-collection'. They run 'sh stats.sh', which outputs a table of metrics. The table has columns: timestamp, cpu_perc, mem_usage, and mem_perc. The data shows metrics for 2026-01-02 from 12:17:58 to 12:18:22. The terminal is then suspended with '^Z' and 'zsh: suspended sh stats.sh'.

```
lolitha@Lolithas-MacBook-Air stats-metric-collection % sh stats.sh
timestamp,cpu_perc,mem_usage,mem_perc
2026-01-02 12:17:58,5.40%,804MiB / 5.787GiB,13.57%
2026-01-02 12:18:01,4.21%,806.2MiB / 5.787GiB,13.61%
2026-01-02 12:18:03,4.47%,806.4MiB / 5.787GiB,13.61%
2026-01-02 12:18:06,7.89%,806.7MiB / 5.787GiB,13.61%
2026-01-02 12:18:09,125.39%,812.8MiB / 5.787GiB,13.72%
2026-01-02 12:18:11,7.61%,815.5MiB / 5.787GiB,13.76%
2026-01-02 12:18:14,6.55%,815.4MiB / 5.787GiB,13.76%
2026-01-02 12:18:16,5.71%,815.6MiB / 5.787GiB,13.76%
2026-01-02 12:18:19,5.32%,815.7MiB / 5.787GiB,13.77%
2026-01-02 12:18:22,3.50%,815.7MiB / 5.787GiB,13.77%
^Z
zsh: suspended  sh stats.sh
lolitha@Lolithas-MacBook-Air stats-metric-collection %
```



Memory usage stays mostly steady across the run. This suggests there is no memory pressure. CPU usage shows short spikes at a few time points. These spikes likely come from bursty input and window processing. Overall changes in performance are driven more by CPU bursts than by memory growth.

Task 5

Application of Big Data Technologies in the Development, Verification and Deployment of Large Language Models (LLMs)

Big Data is field of handling datasets with enormous volume, high velocity, higher variety, higher veracity and much of a potential value. The data driven decision making, which is a crucial process of almost all the industries of fields in the current world highly relies on the information retrieved and generated by processing these Big Data. This process requires key techniques and technologies such as distributed systems and parallel processing frameworks to play with the massive and diverse datasets. Large Language Models is a recently advanced Artificial Intelligence technology which can understand, generating and interacting with human languages using NLP techniques. These LLMs are called “Large” because they are trained on massive datasets with billions or even trillions of parameters to learn and understand the different traits of human languages and to generate response in more humanized manner. The integration of Big Data and Large Language Models comes where the Big Data principles and technologies become the fundamental to all the stages of LLMs, since the primary challenges in Development, Verification and Deployment of LLMs are to manage and handle the massive text corpora used in there.

The Development and Training phase of LLMs involves ingestion of petabyte-scale and diverse text datasets into the LLM pipeline. Parallel data acquisition, preprocessing and distributed training optimization can be distinguished as the two critical areas in this development and training phase. The Big Data frameworks for Petabyte-Scale data pipelines comes to the play to handle massively parallel data acquisition and preprocessing. The distributed platforms like Apache Spark, a tool for large scale data Extract-Transform-Load (ETL) process, Ray Data & Dask, a framework offers parallelization and execution of larger-than-memory and the specialized toolkits like Data Prep Kit (DPK), which is an open-source toolkit to streamline the complex process of quality assessment and transformation of

natural language and code data for LLMs can be identified as the major frameworks for handle large scale data processes[1]. The above platforms ensure the advanced data hygiene by handling the scale, complexity and fuzzy deduplication (identifying and merging approximately same records). The distributed training optimization, which implies overcoming the hardware limitations for during the development process is acquired by using the frameworks such as DeepSeed, which reduces the memory consumption with the Big Data techniques such as Zero Redundancy Optimizer (ZeRO) which partitions the model states across distributed nodes in GPU. This process makes the development and training much efficient by ensuring that each state of the model is stored only once[2].

In the Verification phase of the LLMs the Big Data methodologies are applied to test the LLMs against its biasness, potential failures and other shortcomings at a large scale. This phase uses Big Data applications for three distinct purposes: standardize comparisons for quantitative model performance, large-scale bias, fairness and robustness assessment and synthesizing data for testing and development. The verification process of LLMs needs standardized, large datasets and high-throughput frameworks. Standardized Benchmarking is one of the main tasks to get standardized tasks and datasets to quantitatively assess the reasoning, language understanding and factual accuracy of LLM models [3]. Super General Language Understanding Evaluation (SuperGLUE) is one of the most popular benchmark designed for performance evaluation in Natural Language Understanding (NLU) in LLMs. Also, the use of Big Data techniques to continuously refresh, govern and create unpolluted dataset to mitigate the data contamination issue can be seen as one of the key applications [4] . The bias evaluation process of the LLMs has now shifted to internal bias mitigation strategies, typically involve computationally intensive processes such as detecting and normalizing sensitive attributes within the model [5]. These strategies require distributed systems to handle the scaled analysis of the model's internal data flow.

The final phase of LLM lifecycle is the Deployment phase, typically known as LLM Deployment and Inference (LLMOps). This phase gives the Big Data systems challenge of

managing the high cost and low-latency demands of large models. The optimization of the GPU resources is a crucial need to manage the cost at deployment and there comes the use of various Big Data techniques. One of the key techniques used for this purpose is PagedAttention which uses a core Big Data resource optimization principle of effectively eliminating fragmentation to boost the throughput significantly [6]. This technique mainly applies the distributed memory management theory to LLM's Key-Value cache to partition the cache to allocate the GPU only when needed. Also, GPU utilization is maximized by processing incoming requests in parallel using Continuous Batching techniques which uses the parallel processing principles in Big Data. Furthermore, for the multi-node performance of the LLMs, these systems use Distributing Serving Architectures which share cache as shared to handle highly volatile distributed data resources at the deployment phase. However, still there are some ethical and big data challenges such as, LLMs memorizing personal information and the black box nature of LLM responses which are to be handled in deployment phase.

The above comprehensive summary of the applications on Big Data technologies and techniques throughout all the phases in LLM life cycle shows that, Big Data is not just a supplement to LLMs, but those theories are the foundation for the existence of the LLMs since they are obviously “Large” in nature and work with enormous amounts of data at each phase. In brief, the development phase relies on the distributed processing frameworks and system optimization techniques to execute advanced data hygiene and to manage massive memory demands, the verification phase transforms into a Big Data problem of utilizing high-throughput distributed systems and execute internal bias mitigation strategies to evaluate the LLMs quantitatively and finally, the deployment phase mitigate the challenges of feasibility by utilizing the distributed systems concepts such as memory virtualization and dynamic scheduling principles. Therefore, it is evident that many of Big Data concepts do a crucial intervention through Development, Verification and Deployment phases of the LLMs.

References:

- [1] D. Wood et al., “Data-Prep-Kit: getting your data ready for LLM application development,” Nov. 13, 2024, arXiv: arXiv:2409.18164. doi: 10.48550/arXiv.2409.18164.
- [2] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He, “ZeRO: Memory Optimizations Toward Training Trillion Parameter Models,” May 13, 2020, arXiv: arXiv:1910.02054. doi: 10.48550/arXiv.1910.02054.
- [3] S. Chen et al., “Recent Advances in Large Language Model Benchmarks against Data Contamination: From Static to Dynamic Evaluation,” Sept. 29, 2025, arXiv: arXiv:2502.17521. doi: 10.48550/arXiv.2502.17521.
- [4] Y. Cheng, Y. Chang, and Y. Wu, “A Survey on Data Contamination for Large Language Models,” June 05, 2025, arXiv: arXiv:2502.14425. doi: 10.48550/arXiv.2502.14425.
- [5] A. Abhishek, L. Erickson, and T. Bandopadhyay, “BEATS: Bias Evaluation and Assessment Test Suite for Large Language Models,” Mar. 31, 2025, arXiv: arXiv:2503.24310. doi: 10.48550/arXiv.2503.24310.
- [6] W. Kwon et al., “Efficient Memory Management for Large Language Model Serving with PagedAttention,” Sept. 12, 2023, arXiv: arXiv:2309.06180. doi: 10.48550/arXiv.2309.06180.