VIET NAM NATIONAL UNIVERSITY HO CHI MINH CITY
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



**SEMESTER 241 - CO3071**

---

**LAB 5 REPORT**

**DISTRIBUTED SYSTEMS**

---

**Instructor**: **PROF. Nguyễn Mạnh Thìn**

| | |
|---|---|
| Tạ Gia Khang | - 2152642 |
| Nguyễn Quang Thiện | - 2152994 |

HO CHI MINH CITY,  DECEMBER 2024

# Contents

# 1 Homework

## 1.1 Configuration

Using Docker and generating a docker-compose.yml file to create 3 brokers for **Kafka Fog**, **Spark Master**, **2 Spark Workers**, **a Hadoop Namenode**, and **3 Datanodes**.

```yaml
version: '2'
services:
# kafka-fog
  kafka1:
    image: apache/kafka:3.7.0
    hostname: kafka1
    container_name: kafka1
    ports:
      - 29092:9092
    environment:
      KAFKA_NODE_ID: 1
      KAFKA_PROCESS_ROLES: 'broker,controller'
      KAFKA_LISTENER_SECURITY_PROTOCOL_MAP:
        'CONTROLLER:PLAINTEXT,PLAINTEXT:PLAINTEXT,PLAINTEXT_HOST:PLAINTEXT'
      KAFKA_CONTROLLER_QUORUM_VOTERS: '1@kafka1:9093,2@kafka2:9093,3@kafka3:9093'
      KAFKA_LISTENERS:
        'PLAINTEXT://:29092,CONTROLLER://:9093,PLAINTEXT_HOST://:9092'
      KAFKA_INTER_BROKER_LISTENER_NAME: 'PLAINTEXT'
      KAFKA_ADVERTISED_LISTENERS:
        PLAINTEXT://kafka1:29092,PLAINTEXT_HOST://localhost:29092
      KAFKA_CONTROLLER_LISTENER_NAMES: 'CONTROLLER'
      CLUSTER_ID: '2UKPPoqNTPezeassrAogQw'
      KAFKA_OFFSETS_TOPIC_NUM_PARTITIONS: 1
      KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
      KAFKA_GROUP_INITIAL_REBALANCE_DELAY_MS: 0
      KAFKA_TRANSACTION_STATE_LOG_MIN_ISR: 1
      KAFKA_TRANSACTION_STATE_LOG_REPLICATION_FACTOR: 1
      KAFKA_LOG_DIRS: '/tmp/kraft-combined-logs'
    networks:
      - lab5_ds_kafka_network

  kafka2:
    image: apache/kafka:3.7.0
    hostname: kafka2
    container_name: kafka2
    ports:
      - 39092:9092
    environment:
      KAFKA_NODE_ID: 2
      KAFKA_PROCESS_ROLES: 'broker,controller'
```

```yaml
38        KAFKA_LISTENER_SECURITY_PROTOCOL_MAP:
          'CONTROLLER:PLAINTEXT,PLAINTEXT:PLAINTEXT,PLAINTEXT_HOST:PLAINTEXT'
39        KAFKA_CONTROLLER_QUORUM_VOTERS: '1@kafka1:9093,2@kafka2:9093,3@kafka3:9093'
40        KAFKA_LISTENERS:
          'PLAINTEXT://:39092,CONTROLLER://:9093,PLAINTEXT_HOST://:9092'
41        KAFKA_INTER_BROKER_LISTENER_NAME: 'PLAINTEXT'
42        KAFKA_ADVERTISED_LISTENERS:
          PLAINTEXT://kafka2:39092,PLAINTEXT_HOST://localhost:39092
43        KAFKA_CONTROLLER_LISTENER_NAMES: 'CONTROLLER'
44        CLUSTER_ID: '2UKPPoqNTPezeassrAogQw'
45        KAFKA_OFFSETS_TOPIC_NUM_PARTITIONS: 1
46        KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
47        KAFKA_GROUP_INITIAL_REBALANCE_DELAY_MS: 0
48        KAFKA_TRANSACTION_STATE_LOG_MIN_ISR: 1
49        KAFKA_TRANSACTION_STATE_LOG_REPLICATION_FACTOR: 1
50        KAFKA_LOG_DIRS: '/tmp/kraft-combined-logs'
51      networks:
52        - lab5_ds_kafka_network
53
54    kafka3:
55      image: apache/kafka:3.7.0
56      hostname: kafka3
57      container_name: kafka3
58      ports:
59        - 49092:9092
60      environment:
61        KAFKA_NODE_ID: 3
62        KAFKA_PROCESS_ROLES: 'broker,controller'
63        KAFKA_LISTENER_SECURITY_PROTOCOL_MAP:
          'CONTROLLER:PLAINTEXT,PLAINTEXT:PLAINTEXT,PLAINTEXT_HOST:PLAINTEXT'
64        KAFKA_CONTROLLER_QUORUM_VOTERS: '1@kafka1:9093,2@kafka2:9093,3@kafka3:9093'
65        KAFKA_LISTENERS:
          'PLAINTEXT://:49092,CONTROLLER://:9093,PLAINTEXT_HOST://:9092'
66        KAFKA_INTER_BROKER_LISTENER_NAME: 'PLAINTEXT'
67        KAFKA_ADVERTISED_LISTENERS:
          PLAINTEXT://kafka3:49092,PLAINTEXT_HOST://localhost:49092
68        KAFKA_CONTROLLER_LISTENER_NAMES: 'CONTROLLER'
69        CLUSTER_ID: '2UKPPoqNTPezeassrAogQw'
70        KAFKA_OFFSETS_TOPIC_NUM_PARTITIONS: 1
71        KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
72        KAFKA_GROUP_INITIAL_REBALANCE_DELAY_MS: 0
73        KAFKA_TRANSACTION_STATE_LOG_MIN_ISR: 1
74        KAFKA_TRANSACTION_STATE_LOG_REPLICATION_FACTOR: 1
75        KAFKA_LOG_DIRS: '/tmp/kraft-combined-logs'
76      networks:
77        - lab5_ds_kafka_network
78
```

```yaml
79   kafka-ui:
80     container_name: kafka-ui
81     image: provectuslabs/kafka-ui:latest
82     ports:
83       - "8080:8080"
84     environment:
85       DYNAMIC_CONFIG_ENABLED: 'true'
86     depends_on:
87       - kafka1
88       - kafka2
89       - kafka3
90     networks:
91       - lab5_ds_kafka_network
92
93   spark-master:
94     image: bitnami/spark:latest
95     container_name: spark-master
96     environment:
97       - SPARK_MODE=master
98       - SPARK_RPC_AUTHENTICATION_ENABLED=no
99       - SPARK_RPC_ENCRYPTION_ENABLED=no
100      - SPARK_LOCAL_STORAGE_ENCRYPTION_ENABLED=no
101      - SPARK_SSL_ENABLED=no
102    ports:
103      - "7077:7077" # Spark Master port
104    networks:
105      - lab5_ds_kafka_network
106    volumes:
107      - ./pyspark:/opt/spark/python
108
109  spark-worker1:
110    image: bitnami/spark:latest
111    container_name: spark-worker1
112    environment:
113      - SPARK_MODE=worker
114      - SPARK_MASTER_URL=spark://spark-master:7077
115    depends_on:
116      - spark-master
117    networks:
118      - lab5_ds_kafka_network
119
120  spark-worker2:
121    image: bitnami/spark:latest
122    container_name: spark-worker2
123    environment:
124      - SPARK_MODE=worker
125      - SPARK_MASTER_URL=spark://spark-master:7077
```

```yaml
126       depends_on:
127         - spark-master
128       networks:
129         - lab5_ds_kafka_network
130
131   namenode:
132     image: bde2020/hadoop-namenode:latest
133     container_name: namenode
134     environment:
135       - CLUSTER_NAME=hadoop_cluster
136       - CORE_CONF_fs_defaultFS=hdfs://namenode:9000
137     ports:
138       - "9000:9000" # HDFS namenode RPC port
139       - "9870:9870" # HDFS namenode Web UI
140     networks:
141       - lab5_ds_kafka_network
142     volumes:
143       - ./namenode:/hadoop/dfs/name
144
145   datanode:
146     image: bde2020/hadoop-datanode:latest
147     container_name: datanode
148     environment:
149       - CORE_CONF_fs_defaultFS=hdfs://namenode:9000
150     depends_on:
151       - namenode
152     networks:
153       - lab5_ds_kafka_network
154     volumes:
155       - ./datanode:/hadoop/dfs/data
156
157   datanode2:
158     image: bde2020/hadoop-datanode:latest
159     container_name: datanode2
160     environment:
161       - CORE_CONF_fs_defaultFS=hdfs://namenode:9000
162     depends_on:
163       - namenode
164     networks:
165       - lab5_ds_kafka_network
166     volumes:
167       - ./datanode2:/hadoop/dfs/data
168
169   datanode3:
170     image: bde2020/hadoop-datanode:latest
171     container_name: datanode3
172     environment:
```

```
173        - CORE_CONF_fs_defaultFS=hdfs://namenode:9000
174      depends_on:
175        - namenode
176      networks:
177        - lab5_ds_kafka_network
178      volumes:
179        - ./datanode3:/hadoop/dfs/data
180
181  networks:
182    lab5_ds_kafka_network:
183      external: true
184
185  volumes:
186    # Defines volumes for persistent storage of Namenode and Datanode data
187    namenode:
188    datanode:
189    datanode2:
190    datanode3:
```

After that, create the network **lab5_ds_kafka_network** by the cmd below in Docker:

```
1  docker network create lab5_ds_kafka_network
```

## 1.2   Previous LAB's Output

```
1  Check Output:
2    docker exec -it namenode bash -c 'hdfs dfs -ls /hdfs/environment_data'
3    docker exec -it namenode bash -c 'hdfs dfs -cat /hdfs/environment_data
      /*.json'
```
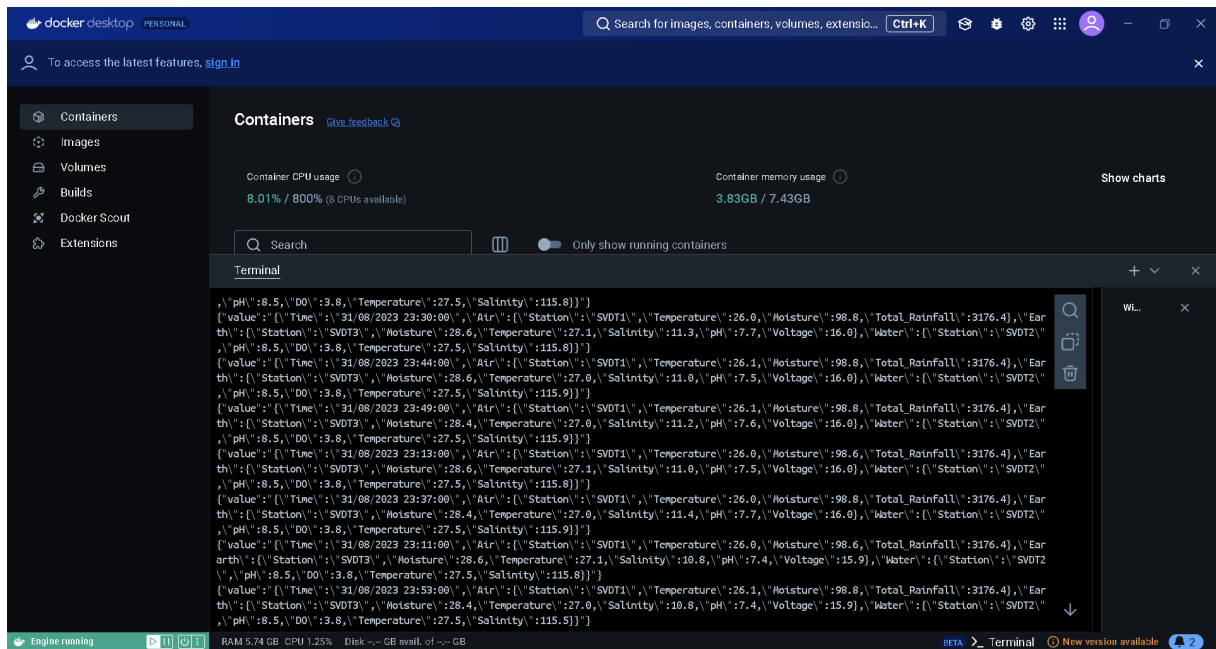


Figure 1: LAB4's Output

## 1.3 Loading Data

This section focuses on initializing the SparkSession and loading the environmental data stored in HDFS for distributed processing. The primary goal is to prepare the data for analysis of mango growth conditions.

- SparkSession serves as the entry point for PySpark applications, enabling distributed processing of large datasets. The session is named *"Distributed Computing for Mango Growth Evaluation"* for easy identification.

- The input data, stored as JSON files which contains information about air temperature, air moisture, and soil pH in HDFS at hdfs://namenode:9000/hdfs/environment_data/*.json.

- A schema is defined using StructType to parse the nested JSON structure accurately. It includes:

    - *Air*: Contains fields for temperature and moisture.
    - *Earth*: Contains the field for soil pH.

- The loaded data is parsed using from_json, and generating structured columns Air_Temperature, Air_Moisture, and Earth_pH, which are really vital for the subsequent analysis.

The parsed and structured dataset enables distributed computation while maintaining scalability and fault tolerance, leveraging the advantages of HDFS and Spark for efficient data handling.

```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, from_json, when, lit, avg, max, min
from pyspark.sql.types import StructType, StructField, DoubleType, StringType

# Initialize SparkSession
spark = SparkSession.builder \
    .appName("Distributed Computing for Mango Growth Evaluation") \
    .getOrCreate()

# HDFS Path
input_path = "hdfs://namenode:9000/hdfs/environment_data/*.json"
output_path = "hdfs://namenode:9000/hdfs/mango_data_distributed.txt"

# Define Schema for Nested JSON
schema = StructType([
    StructField("Air", StructType([
        StructField("Temperature", DoubleType()),
        StructField("Moisture", DoubleType())
    ])),
    StructField("Earth", StructType([
        StructField("pH", DoubleType())
    ]))
])

# Define Ideal Conditions
ideal_conditions = {
```

```python
27        "Air_Temperature": (23.0, 28.0),
28        "Air_Moisture": (60.0, 80.0),
29        "Earth_pH": (5.5, 7.0)
30    }
31
32    def process_data():
33        # Load JSON Data and collect to driver
34        df = spark.read.json(input_path)
35
36        # Parse the "value" field
37        parsed_df = df.withColumn("parsed_value", from_json(col("value"), schema))
38
39        # Extract the parsed fields
40        data_df = parsed_df.select(
41            col("parsed_value.Air.Temperature").alias("Air_Temperature"),
42            col("parsed_value.Air.Moisture").alias("Air_Moisture"),
43            col("parsed_value.Earth.pH").alias("Earth_pH")
44        )
```

## 1.4   Bring data to compute (Centralized Computing)

This section explains the centralized computation approach, where the data is first collected from the distributed HDFS system and processed locally.

- **Data Collection:** All records from the distributed dataset are gathered to the driver node using the `collect()` method. This transforms the distributed data into a local collection, enabling computation to occur in a centralized environment.

- **Identifying Out-of-Ideal Conditions:** The script evaluates whether environmental parameters (`Air_Temperature`, `Air_Moisture`, `Earth_pH`) fall outside their ideal ranges:

  - `Air_Temperature`: Between 23°C and 28°C.
  - `Air_Moisture`: Between 60% and 80%.
  - `Earth_pH`: Between 5.5 and 7.0.

  Each parameter is flagged with a `1` if it is out of range and `0` otherwise.

- **Real-Time Monitoring:** Using `show()` function outputs the DataFrame to monitor conditions in real-time, providing insights into the current state of the environment.

- **Time Analysis:** The total number of rows (representing minutes) where one or more parameters are out of range is counted. This is converted into hours and minutes for ease of understanding.

- **Statistical Analysis:** Using `agg()` function, the maximum and minimum values for each parameter are computed to assess the range and extremes of environmental conditions, offering valuable insights into the dataset's variability.

This method highlights the trade-off of centralized computing: while simple to implement, it may not scale efficiently for large datasets due to data transfer overheads.

```python
1    # Collect all data to driver node
2    collected_data = data_df.collect()
3
4    # Convert collected data to a new local DataFrame
5    local_df = spark.createDataFrame(collected_data)
6
7    # Calculate Conditions Out of Ideal
8    out_of_ideal_df = local_df.withColumn(
9        "Air_Temperature_Out",
10       when((col("Air_Temperature") <
         lit(ideal_conditions["Air_Temperature"][0])) |
11           (col("Air_Temperature") >
             lit(ideal_conditions["Air_Temperature"][1])), 1).otherwise(0)
12   ).withColumn(
13       "Air_Moisture_Out",
14       when((col("Air_Moisture") < lit(ideal_conditions["Air_Moisture"][0])) |
15           (col("Air_Moisture") > lit(ideal_conditions["Air_Moisture"][1])),
             1).otherwise(0)
16   ).withColumn(
17       "Earth_pH_Out",
```

```python
18          when((col("Earth_pH") < lit(ideal_conditions["Earth_pH"][0])) |
19              (col("Earth_pH") > lit(ideal_conditions["Earth_pH"][1])),
                1).otherwise(0)
20      )
21
22      # Calculate time out of ideal conditions locally
23      out_of_ideal_count = out_of_ideal_df.filter(
24          (col("Air_Temperature_Out") + col("Air_Moisture_Out") +
            col("Earth_pH_Out")) >= 1
25      ).count()
26
27      # Convert from minutes to hours, minutes
28      hours = out_of_ideal_count // 60
29      minutes = out_of_ideal_count % 60
30
31      # Calculate statistics locally
32      stats = local_df.agg(
33          max("Air_Temperature").alias("Max_Air_Temperature"),
34          min("Air_Temperature").alias("Min_Air_Temperature"),
35          max("Air_Moisture").alias("Max_Air_Moisture"),
36          min("Air_Moisture").alias("Min_Air_Moisture"),
37          max("Earth_pH").alias("Max_Earth_pH"),
38          min("Earth_pH").alias("Min_Earth_pH")
39      ).collect()[0]
```

## 1.5 Bring compute to data (Distributed Computing)

This section demonstrates the distributed computing approach, where the computation is brought to the data stored in HDFS, reducing data transfer overhead and leveraging distributed resources.

- **Identifying Out-of-Ideal Conditions:** Similar to the centralized approach, new columns are added to evaluate whether `Air_Temperature`, `Air_Moisture`, and `Earth_pH` fall outside their respective ideal ranges. This computation occurs in parallel across the distributed nodes, enhancing efficiency.

  - `Air_Temperature`: Ideal range is between 23°C and 28°C.
  - `Air_Moisture`: Ideal range is between 60% and 80%.
  - `Earth_pH`: Ideal range is between 5.5 and 7.0.

- **Real-Time Monitoring:** `show()` function is also utilized to display the DataFrame, allowing real-time observation of the conditions and offering insights into the current environmental state.

- **Time Analysis:** The script calculates the total time (in minutes) where one or more parameters are outside the ideal range. This is further converted to hours and minutes for better interpretability.

- **Statistical Summary:** The maximum and minimum values for each parameter are computed using Spark's `agg()` function. These calculations are distributed across the cluster, reducing processing time while handling large-scale data.

This approach effectively utilizes distributed computing capabilities to analyze environmental data at scale, minimizing the cost of moving large datasets.

```python
# Calculate Conditions Out of Ideal
out_of_ideal_df = data_df.withColumn(
    "Air_Temperature_Out",
    when((col("Air_Temperature") < lit(ideal_conditions["Air_Temperature"][0])) |
        (col("Air_Temperature") > lit(ideal_conditions["Air_Temperature"][1])),
        1).otherwise(0)
).withColumn(
    "Air_Moisture_Out",
    when((col("Air_Moisture") < lit(ideal_conditions["Air_Moisture"][0])) |
        (col("Air_Moisture") > lit(ideal_conditions["Air_Moisture"][1])),
        1).otherwise(0)
).withColumn(
    "Earth_pH_Out",
    when((col("Earth_pH") < lit(ideal_conditions["Earth_pH"][0])) |
        (col("Earth_pH") > lit(ideal_conditions["Earth_pH"][1])), 1).otherwise(0)
)

# Show the DataFrame
out_of_ideal_df.show()

# Calculate Average Time Out of Ideal Conditions
avg_out_of_ideal = out_of_ideal_df.filter((col("Air_Temperature_Out") +
col("Air_Moisture_Out") + col("Earth_pH_Out")) >= 1).count()

```

```python
22  # Convert from minutes to hours, minutes, seconds
23  avg_out_of_ideal_hours = avg_out_of_ideal // 60
24  avg_out_of_ideal_minutes = avg_out_of_ideal % 60
25
26  # Calculate Max and Min for Each Field
27  stats_df = data_df.agg(
28      max("Air_Temperature").alias("Max_Air_Temperature"),
29      min("Air_Temperature").alias("Min_Air_Temperature"),
30      max("Air_Moisture").alias("Max_Air_Moisture"),
31      min("Air_Moisture").alias("Min_Air_Moisture"),
32      max("Earth_pH").alias("Max_Earth_pH"),
33      min("Earth_pH").alias("Min_Earth_pH")
34  ).collect()[0]
```

## 1.6 Store Output

The final step involves saving the processed results back to HDFS in a structured format for future use and further analysis.

- **Formatting Results:** The computed statistics and time metrics are formatted into a readable string that summarizes the analysis.

- **Storing Results:** The function `saveAsTextFile()` writes the formatted results as a single text file to the specified HDFS directory. This ensures the outputs are accessible and reusable for downstream processes.

- **Finalization:** SparkSession is terminated to free up resources, signifying the completion of the computation process.

This step ensures that all insights derived from the distributed computing process are preserved in a distributed storage system, enabling easy sharing and accessibility.

```python
    # Format results
    results = (
        f"Average Time Out of Ideal: {hours} hours {minutes} minutes\n"
        f"Max Air Temperature: {stats['Max_Air_Temperature']}°C\n"
        f"Min Air Temperature: {stats['Min_Air_Temperature']}°C\n"
        f"Max Air Moisture: {stats['Max_Air_Moisture']}%\n"
        f"Min Air Moisture: {stats['Min_Air_Moisture']}%\n"
        f"Max Earth pH: {stats['Max_Earth_pH']}\n"
        f"Min Earth pH: {stats['Min_Earth_pH']}\n"
    )

    # Write results to HDFS as a single text file
    spark.sparkContext.parallelize([results]).saveAsTextFile(output_path)

# Process the data
process_data()

# Stop SparkSession
spark.stop()
```

## 1.7 Demonstration

### 1.7.1 Centralized Computing Method

```
1 Push code to Spark:
2   docker cp data_to_compute.py spark-master:/opt/spark/python/
      data_to_compute.py
3 Run code:
4   docker exec -it spark-master bash
5   spark-submit \
6     --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.3 \
7     --master spark://spark-master:7077 \
8   /opt/spark/python/data_to_compute.py
```
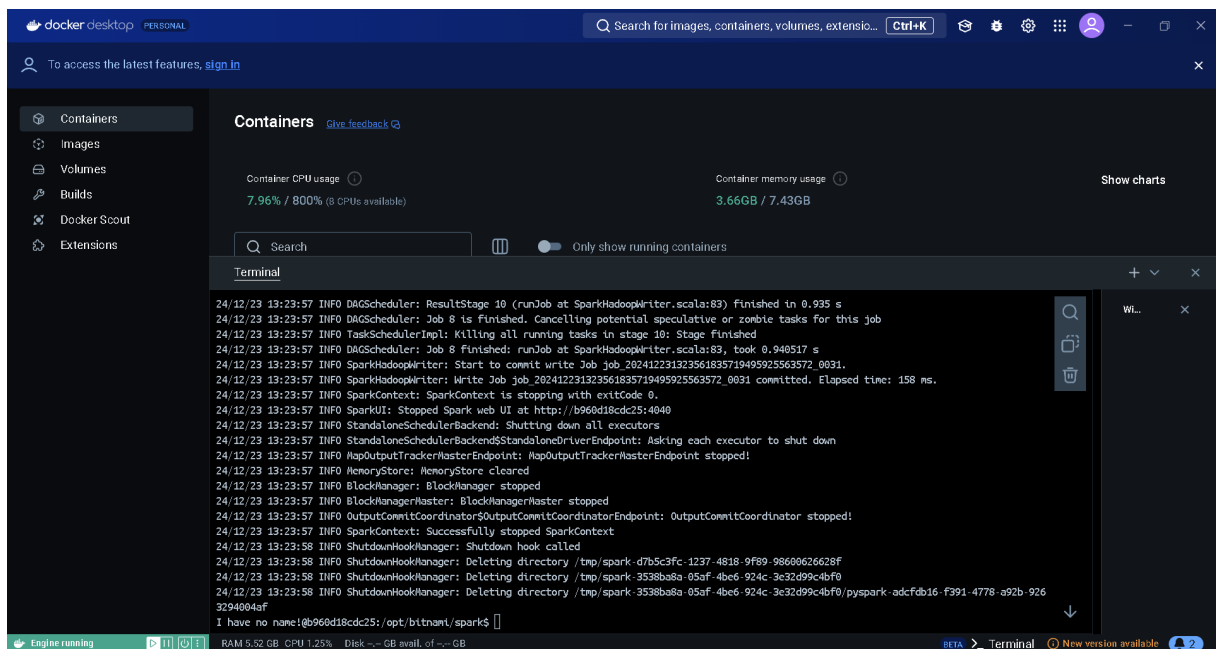


Figure 2: Finish Centralized Computing

```
1 Check Output:
2   docker exec -it namenode bash -c 'hdfs dfs -cat /hdfs/
      mango_data_centralized.txt/part-*'
```
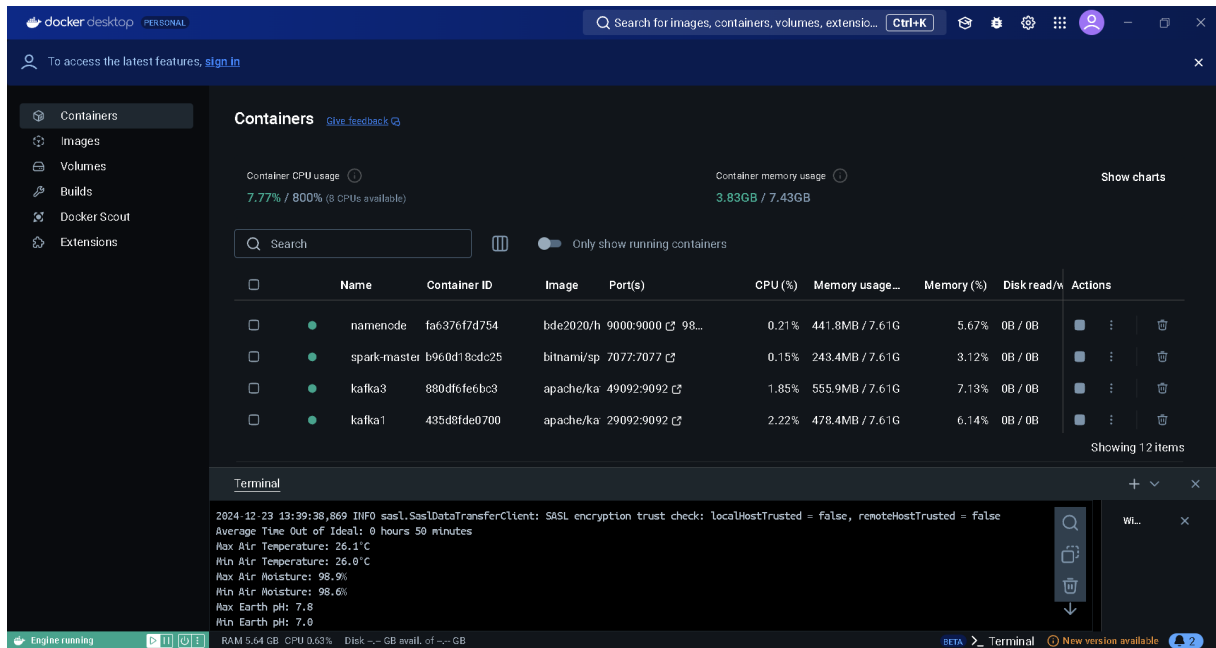
Figure 3: Output Centralized Computing

### 1.7.2 Distributed Computing Method

```
Push code to Spark:
  docker cp compute_to_data.py spark-master:/opt/spark/python/
    compute_to_data.py
Run code:
  docker exec -it spark-master bash
  spark-submit \
    --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.3 \
    --master spark://spark-master:7077 \
  /opt/spark/python/compute_to_data.py
```
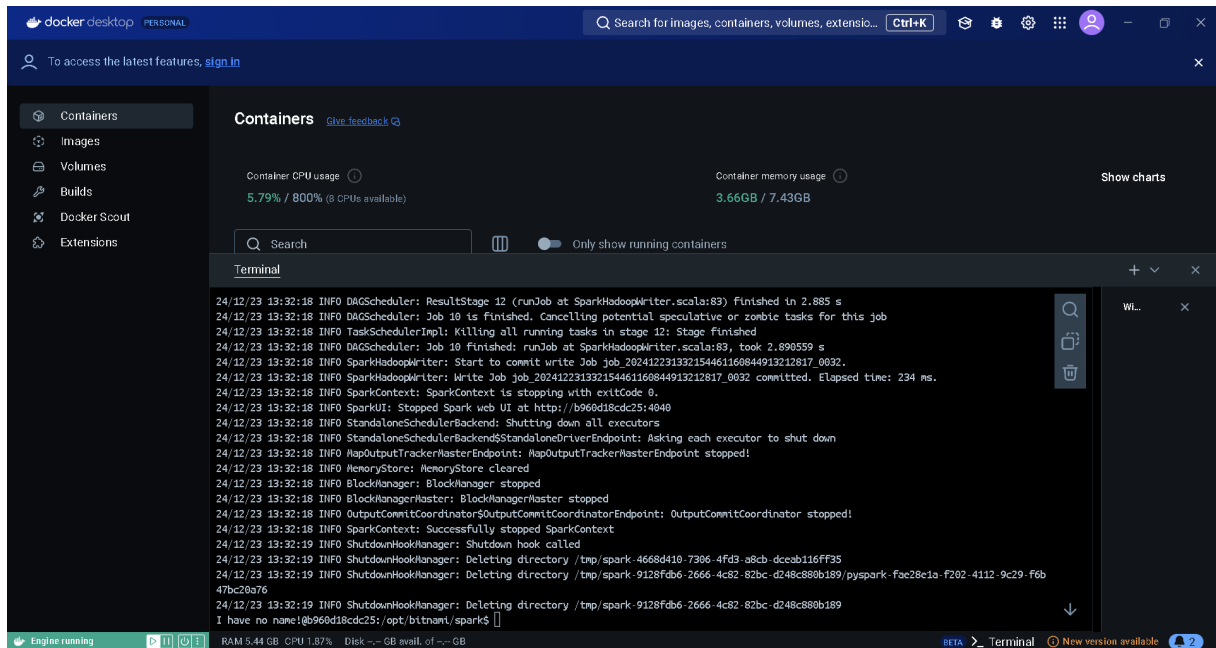
Figure 4: Finish Distributed Computing

```
Check Output:
  docker exec -it namenode bash -c 'hdfs dfs -cat /hdfs/
    mango_data_distributed.txt/part-*'
```
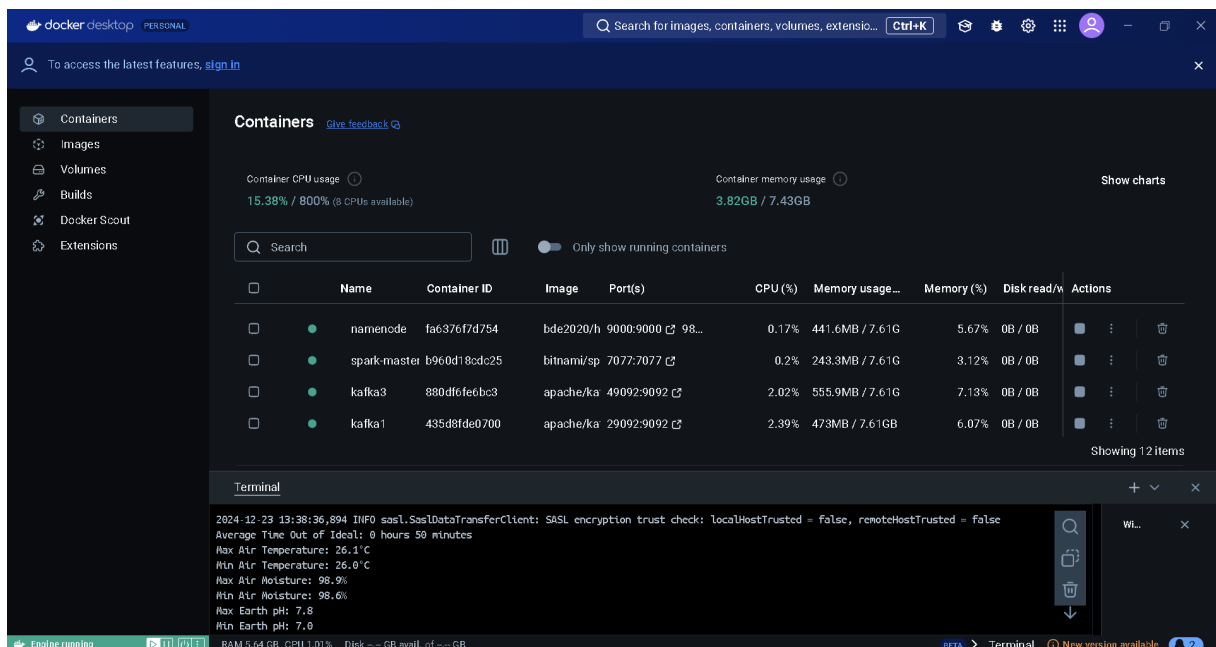


Figure 5: Output Distributed Computing

Observing the two output images generated by the two methods reveals that the results of both approaches are identical.