

REPORT

# Multi-language Engine: Apache Spark

Big Data Analytics and Reasoning

WRITTEN BY Kreuzer Florian, Matousch David

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Dataset Description</b>	<b>2</b>
2.1	Example Request . . . . .	2
<b>3</b>	<b>Implementation</b>	<b>3</b>
3.1	Batch processing . . . . .	3
3.2	Stream processing . . . . .	4
3.2.1	Data generator . . . . .	4
3.2.2	Data processor . . . . .	5
<b>4</b>	<b>Results and Insights</b>	<b>5</b>
<b>5</b>	<b>Challenges and Solutions</b>	<b>5</b>
5.1	Setup . . . . .	5
5.2	Batch processing . . . . .	5
5.3	Stream processing . . . . .	6
<b>6</b>	<b>Deployment Details</b>	<b>6</b>
6.1	How to run . . . . .	6
6.2	Deployment Implementation . . . . .	6
<b>7</b>	<b>Conclusion</b>	<b>7</b>

## 1 Introduction

Efficient processing of large-scale datasets has become a critical requirement across industries. Apache Spark addresses this need by offering a distributed computing framework designed for speed and scalability. Unlike Hadoop, which relies heavily on disk storage and HDFS for computations, Spark operates primarily in memory, significantly improving performance for iterative tasks.

Initially, Spark introduced Resilient Distributed Datasets (RDDs) as its core abstraction for distributed data processing. However, since Spark 2.0, RDDs have been largely replaced by the Dataset API, which combines the benefits of RDDs with optimizations from Spark's Catalyst query optimizer. This shift reflects a move toward a more structured and efficient approach to data processing.

The Spark stack includes Spark Core, which provides essential functionalities such as task scheduling, memory management, and fault recovery. Built on top of Spark Core are higher-level APIs like DataFrames and Structured Streaming, enabling developers to process structured and unstructured data efficiently. (see Figure 1)

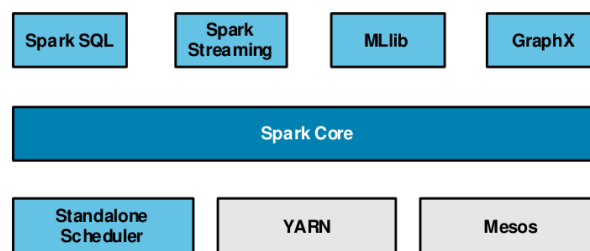


Figure 1: Spark stack and widely used higher-level APIs

## 2 Dataset Description

The dataset used in this project is sourced from the *OpenWeatherMap API*, a comprehensive platform for accessing weather data. This API provides current and forecasted weather information, along with historical data spanning over 45 years. The key features of the dataset include: [1]

- **Current Weather and Forecasts:**
  - Minute-by-minute forecasts for the next hour.
  - Hourly forecasts for up to 48 hours.
  - Daily forecasts for up to 8 days.
  - Government-issued weather alerts.
- **Historical Weather Data:**
  - Daily aggregated data for the past 45+ years.
- **API Access:** Users can make up to 1,000 API calls per day for free.

### 2.1 Example Request

To retrieve weather data, an example API call structure is as follows:

```
api.openweathermap.org/data/3.0/onecall?lat={lat}&lon={lon}&exclude={part}&appid={API key}
```

- `lat` and `lon` specify the latitude and longitude of the location.
- `exclude` filters specific data segments, such as current weather or minute forecasts.
- `API key` is the unique key required to authenticate the user.

## 3 Implementation

### 3.1 Batch processing

We used *DataFrames* for our batch processing demo. As seen in Listing 1 the fetching of the weather data involves no spark components.

Listing 1: Fetch Weather Data

```

1 def fetch_weather_data(cities):
2     """Fetch weather data from OpenWeatherMap API for a list of cities."""
3     weather_data = []
4     for city in cities:
5         try:
6             response = requests.get(BASE_URL, params={"q": city, "appid": API_KEY, "units":
7                 "metric"})
8             if response.status_code == 200:
9                 data = response.json()
10                weather_data.append({
11                    "city": data["name"],
12                    "temperature": data["main"]["temp"],
13                    "weather": data["weather"][0]["main"],
14                    "humidity": data["main"]["humidity"]
15                })
16            else:
17                print(f"Failed to fetch data for {city}: {response.status_code}")
18        except Exception as e:
19            print(f"Error fetching data for {city}: {e}")
20    return weather_data

```

In Listing 2 the weatherdata list is converted into a spark DataFrame.

Listing 2: Convert to DataFrame

```

1 spark = SparkSession.builder \
2     .appName("Batch Weather Analysis") \
3     .getOrCreate()
4
5 weather_data = fetch_weather_data(CITIES)
6
7 weather_df = spark.createDataFrame(weather_data)
8
9 print("Raw Weather Data:")
10 weather_df.show()

```

Now we can leverage spark operations to analyze the data. In Listing 3 some examples are shown.

Listing 3: Spark Operations on DataFrame

```

1 # Average Temperatures per City
2 avg_temp_df = weather_df.groupBy("city").agg(
3     avg("temperature").alias("avg_temperature")
4 )
5 print("Average Temperature per City:")
6 avg_temp_df.show()
7
8 # Identify cities with extreme weather conditions
9 extreme_weather_df = weather_df.withColumn(
10     "extreme_weather",
11     when((col("temperature") > 35) | (col("temperature") < 0), lit("Yes")).otherwise(lit("No"))
12 )
13 print("Cities with Extreme Weather:")
14 extreme_weather_df.filter(col("extreme_weather") == "Yes").show()
15
16 # Group cities by weather condition

```

```

17 weather_group_df = weather_df.groupBy("weather").count()
18 print("Cities Grouped by Weather Condition:")
19 weather_group_df.show()

```

Listing 4: Spark Operations on DataFrame Output

```

1 Average Temperature per City:
2 +-----+-----+
3 |          city|avg_temperature|
4 +-----+-----+
5 |      New York|         -1.4|
6 |      London|         6.74|
7 |      Berlin|         0.97|
8 |      Tokyo|         8.48|
9 |      Sydney|        24.22|
10 |      Moscow|        -4.17|
11 |      Mumbai|        28.99|
12 |    Cape Town|        21.67|
13 |Rio de Janeiro|       23.12|
14 |      Beijing|        -3.06|
15 +-----+-----+
16
17 Cities with Extreme Weather:
18 +-----+-----+-----+-----+-----+
19 | city|humidity|temperature|weather|extreme_weather|
20 +-----+-----+-----+-----+-----+
21 |New York|    52|      -1.4| Clear|          Yes|
22 | Moscow|    74|      -4.17|  Snow|          Yes|
23 | Beijing|    19|      -3.06| Clear|          Yes|
24 +-----+-----+-----+-----+-----+
25
26 +-----+-----+
27 |weather|count|
28 +-----+-----+
29 |  Clear|    4|
30 | Clouds|    3|
31 |   Rain|    1|
32 |   Snow|    1|
33 |  Smoke|    1|
34 +-----+-----+

```

## 3.2 Stream processing

The stream processing application consists of two parts. A data generator and a data processor.

### 3.2.1 Data generator

Weather data gets fetched from a random city every 2 seconds. This data is then stored in a new json file.

Listing 5: Data generator

```

1 def simulate_streaming_data():
2     # remove all old files
3     shutil.rmtree(OUTPUT_DIR)
4     os.makedirs(OUTPUT_DIR)
5
6     count = 0
7     while True:
8         city = random.choice(CITIES)
9         weather_data = fetch_weather(city)
10        if weather_data:
11            # Generate a unique filename using a counter

```

```

12     filename = os.path.join(OUTPUT_DIR, f"weather_{count}.json")
13     with open(filename, "w") as f:
14         json.dump(weather_data, f)
15         print(f"Generated: {filename}")
16         count += 1
17         time.sleep(2)

```

This function is running in a separate thread.

Listing 6: Start thread

```

1 threading.Thread(target=simulate_streaming_data, daemon=True).start()

```

After that you just need to attach a spark listener on the generation folder.

Listing 7: streaming DataFrame

```

1 streaming_df = spark.readStream \
2     .schema(schema) \
3     .json(OUTPUT_DIR)

```

### 3.2.2 Data processor

Because the *streaming\_df* is already a normal DataFrame we can run any processing we want using the normal transformations. For our test application we chose the following.

Listing 8: processing DataFrame

```

1 processed_df = streaming_df \
2     .withColumn("extreme_weather",
3         when((col("temperature") > 35) | (col("temperature") < 0),
4             lit("Yes")).otherwise(lit("No"))) \
5     .withColumn("is_rainy", when(col("weather") == "Rain", lit("Yes")).otherwise(lit("No")))

```

To get the results on the screen we added a simple output to the console.

Listing 9: streaming output

```

1 query = processed_df.writeStream \
2     .outputMode("append") \
3     .format("console") \
4     .option("truncate", "false") \
5     .start()
6
7 query.awaitTermination()

```

## 4 Results and Insights

## 5 Challenges and Solutions

### 5.1 Setup

Setup of spark with docker was the main challenge we faced. Our first approach was using *Scala*. Here we had issues with version mismatches. Some objects were serialized and then deserialized using a different java version, thus failing during the process. After that we switched to the current version that uses python to transform and process dataframes. Setup here was easier but Java version mismatches plagued us here as well. Writing a dockerfile that would work on Arm64 and x86 also proved to be challenging.

### 5.2 Batch processing

Batch processing was straightforward and using the weather API was easy. We ran into no problems here.

### 5.3 Stream processing

It is clear that streaming is more difficult than just batching. The main problem here was to understand how the data generator has to generate the data for spark to process it correctly. Creating new json files for every new batch was the solution.

## 6 Deployment Details

### 6.1 How to run

The source code is hosted on GitHub. (<https://github.com/matouschdavid/SparkPlayground.git>) All the necessary commands are listed in the README.md file. Depending on the internet connection, the first run can take multiple minutes.

### 6.2 Deployment Implementation

For our solution, we used *docker* as a deployment platform. When starting the *docker-compose* file, the environment variable *PROCESSING* has to be set. It is used to change between the batch/stream processing. The whole deployment is described in a *Dockerfile* and *docker-compose*

Listing 10: Dockerfile

```

1 FROM python:3.10-slim
2
3 # Set environment variables
4 ENV SPARK_VERSION=3.4.1
5 ENV HADOOP_VERSION=3
6 ENV PYSPARK_PYTHON=python3
7 ENV PYSPARK_DRIVER_PYTHON=python3
8
9 # Install java
10 # ... (not relevant)
11
12 # Download and install Spark
13 RUN wget
14     https://archive.apache.org/dist/spark/spark-${SPARK_VERSION}/spark-${SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}
15     && \
16     tar -xvf spark-${SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz && \
17     mv spark-${SPARK_VERSION}-bin-hadoop${HADOOP_VERSION} /opt/spark && \
18     rm spark-${SPARK_VERSION}-bin-hadoop${HADOOP_VERSION}.tgz
19
20 # Set Spark environment variables
21 ENV SPARK_HOME=/opt/spark
22 ENV PATH=$SPARK_HOME/bin:$PATH
23
24 # Install Python dependencies
25 ## ... (not relevant)
26
27 # Set Spark environment variables for the runtime
28 ENV _JAVA_OPTIONS="-Djava.io.tmpdir=/tmp"
29 ENV PYSPARK_SUBMIT_ARGS="--driver-java-options='-Djava.io.tmpdir=/tmp' pyspark-shell"
30
31 # Default command
32 CMD ["python", "main.py"]

```

The base is a docker image with python installed. In order to use spark, a java installation is required. The *pyspark* package some environment configuration, which is also described in the Dockerfile. (See Listing 10)

Listing 11: docker-compose.yml

```

1 version: "3.9"

```

```
2 services:
3   spark:
4     build: .
5     container_name: spark-container
6     volumes:
7       - ./app:/app
8     working_dir: /app
9     command: ["python", "main.py", "${PROCESSING}"]
```

The `docker-compose` file (Listing 11) uses the `Dockerfile` for the *spark* service. It also defines the processing mode.

## 7 Conclusion

We think that Spark is only useful for actual big data applications. The scale of our project is way too small to make use of the powerful Spark engine. Knowing this, the difficult setup makes it one of our last choices when developing something similar in the future. The API is nice, but nothing new since C# Linq and Java streams have the same API for years natively in their language with many sources you can attach.

That being said, now that we have a working setup, it is great to just spin up a docker container that executes a job and stores everything to our server if we ever need it.

## References

- [1] OpenWeather. Openweather api documentation. <https://openweathermap.org/api/one-call-3>, 2025. Accessed: 2025-01-08.