# INSTITUT SUPÉRIEUR D'ÉLECTRONIQUE DE PARIS (ISEP)



## II.2414 ADVANCED DATABASES AND

## BIG DATA

## BIG DATA PROJECT

## (Data Lake and Data Pipeline)

## BUHARI ALIYU

## ABUBAKAR UMAR ELNAFATY

## PROF. THIBAUT DE BROCA

## JUNE 2023

# Synopsis

The main goal of the 'Big Data project' is to create a simple end-to-end data architecture that spans the entire data lifecycle, encompassing data ingestion, transformation, and exposition. Within this project, you have the freedom to choose the specific data you wish to retrieve and determine the desired output that will generate value from the input data.

However, it is essential to adhere to a Datalake architecture, as it plays a crucial role in organizing the data effectively, ensuring a streamlined data pipeline, and facilitating data sharing among stakeholders. By following the Datalake architecture principles, we can achieve a well-structured and efficient data ecosystem

## Theme for the Data

The central focus of this project revolves around gathering movies, weather data, and forecast information from multiple sources. The primary goal is to conduct a comprehensive analysis of the current top-rated, upcoming, and trending movies, considering their popularity, vote counts, release dates, and weather conditions. By incorporating weather forecast data, this analysis aims to provide valuable insights for individuals planning to visit a cinema.

By visualizing the cinema data alongside weather forecasts, stakeholders such as consumers, researchers, statisticians, and data scientists can make more informed decisions about their moviegoing experiences. The integration of weather information allows users to consider factors like temperature, precipitation, and other weather-related conditions that can impact their cinema experience.

For example, if someone prefers to enjoy outdoor activities before or after watching a movie, they can check the weather forecast to determine the best time to visit the cinema. If it's expected to rain heavily, they might opt for an indoor movie theater instead of an outdoor cinema. On the other hand, if the weather forecast indicates clear skies and pleasant temperatures, they may choose an outdoor cinema for a unique and enjoyable experience.

Moreover, the combination of cinema data and weather forecasts enables users to apply filters and preferences based on their specific moviegoing desires. They can filter movies based on genre, release date, popularity, and user ratings, and further refine their choices by considering weather conditions. For instance, individuals who prefer to watch horror movies during rainy days can easily identify suitable options by cross-referencing the cinema data with the weather forecast.

Furthermore, this enhanced visualization empowers individuals to plan their cinema outings more effectively by taking into account both movie-related factors and

weather conditions. It provides a comprehensive understanding of the movie landscape and allows users to make informed decisions, ensuring an enjoyable and tailored cinema experience

**Data Sources**

The retrieved data would be gathered in the data lake as raw data. Data sources include:

1. The Movie Database TMDB API
2. Weather Api

A community-built movie and television database is called The Movie Database (TMDB) API. The documentation for the API is available at https://developer.themoviedb.org/docs/getting-started , and it includes instructions on how to connect to our application and fetch data via the API. We have fetched 3 different dataset from the api, the top rated movies, the upcoming movies and trending movies in the cinema. We saved the data as a json file. While the weather api documentation can be access via https://www.weatherapi.com/docs, we fetched the current weather and forecast weather of Paris



Figure 1 shows how to access the movie database API.

The Weather API serves as a valuable source of raw data for this project, providing real-time weather forecasts (for the next 7 days) and current weather information for Paris. By integrating this data into the project's data lake, stakeholders can explore the relationship between weather conditions and cinema trends, gaining insights into moviegoer behavior and enabling informed decision-making in Paris.

Figure 2: Connection to weather API

## Data Lake Organization

We designed our data lake with a structure that adheres to the following naming convention: Each route will look like this:

/{layer}/{group}/{TableName}/{date}/filename

**Layers** can include: - **raw** which is the first phase of the data and is straight from the source system, the data is prepared for use by the Datalake in the second stage, which is *formatted*, and *usage* is the final stage before

**Group**: we use *movies* to store movies file and *weather* to save weather files

**Table Name**: An item in this folder will always have this name and a similar schema, in raw we use *imdb* for the movies and *weatherapi* for weather.

**FileName**: toprated_movies.json, upcoming_movies.json, trending_movies.json, current_weather.json and forecast.json are the file names for the raw data.

**Date**: should be formatted as YYYYMMDD, for instance, 20230611 (11th June 2022).

**For formatted file:** toprated_movies.snappy.parquet, trending_movies.snappy.parquet, upcoming_movies.snappy.parquet, current_weather.snappy.parquet and forecast.snappy.parquet.

Figure 3: Datalake folder structure

## Data Pipeline

The data pipeline plays a crucial role in transferring data from its source to a designated destination, such as a data warehouse. Throughout this journey, the data undergoes transformations and optimizations, ultimately arriving in a state that allows for analysis and the development of valuable business insights. The data pipeline encompasses the essential steps of aggregating, organizing, and moving data effectively.

## Data Pipeline Architecture

1.  **Data Ingestion**: The source data is obtained through APIs. The data sources and ingestion process are detailed in Figure 1 and Figure 2 in the data source section. The collected data is then transferred to the raw folder in the data lake, where it can be further prepared for analysis. Python programming language is used for data ingestion in this project.

2.  **Data Formatting**: In this step, the data is aggregated, cleansed, and manipulated to bring it to a standardized format suitable for analysis. Tools like Pandas are utilized to format the data like location, latitude and longitude for the weather. Figure 5 and Figure 6 illustrate the formatting process applied to data from different sources. The formatted data is then converted to the Parquet data format and compressed using snappy, which offers better processing and analysis capabilities.

By following these steps, the data is transformed and prepared in a way that optimizes its usability for subsequent analysis and visualization.



Figure 5: Formatting data from imdb with Panda

Figure 6: Formatting weather data with Panda

3. **Data Combination:** Apache Spark is utilized in this step to combine the data from different data sources and performed sql analysis. This combined data becomes the foundation for further analysis. It is indexed and made accessible through the Elasticsearch and Kibana to visualize the data, we use the cloud version of the elasticsearch and kibana



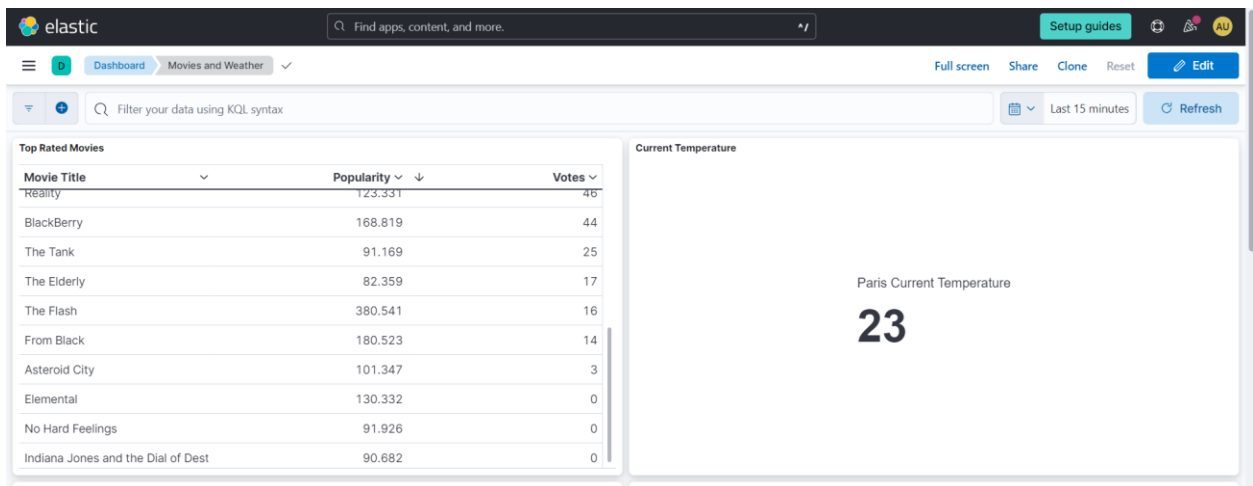Figure 7: Data combination with Apache Spark

Figure 8: Analysis

4. **Indexing**: In this project, the Elasticsearch Python library is utilized to facilitate the indexing process. The data is efficiently indexed and organized, enabling easy retrieval and exploration through Elasticsearch and visualization using Kibana.
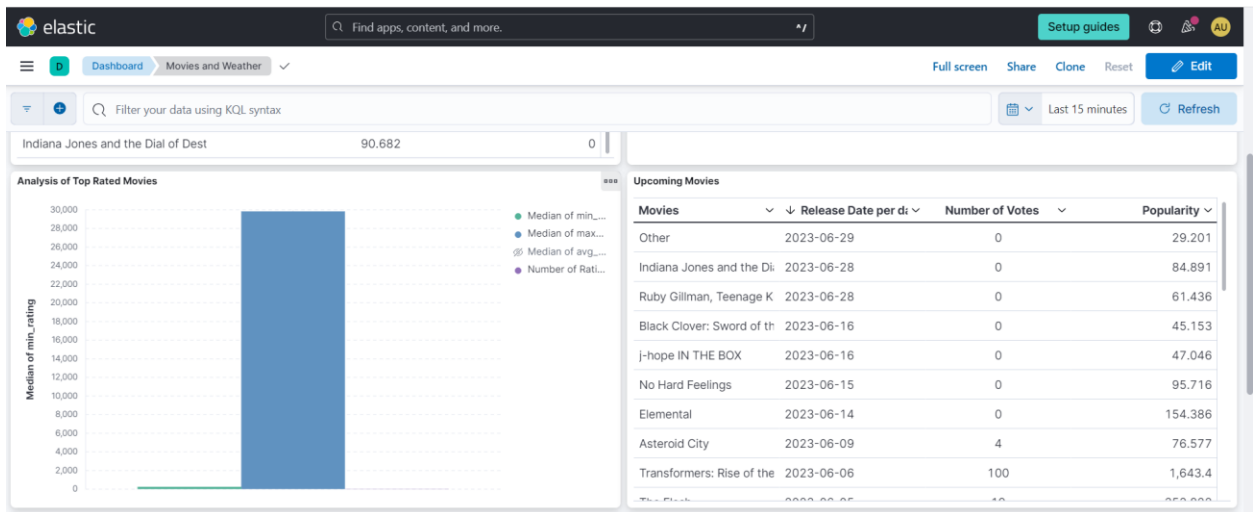
```python
analysis_table = pq.read_table(ANALYSIS_PATH)
df_analysis = analysis_table.to_pandas()
for _, row in df_analysis.iterrows():
    document = row.to_dict()
    es.index(index=analysis_index, body=document)

upcoming_table = pq.read_table(UPCOMING_PATH)
df_upcoming = upcoming_table.to_pandas()
for _, row in df_upcoming.iterrows():
    document = row.to_dict()
    es.index(index=upcoming_index, body=document)

trending_table = pq.read_table(TRENDING_PATH)
df_trending = trending_table.to_pandas()
for _, row in df_trending.iterrows():
    document = row.to_dict()
    es.index(index=trending_index, body=document)

current_table = pq.read_table(CURRENT_PATH)
df_current = current_table.to_pandas()
for _, row in df_current.iterrows():
    document = row.to_dict()
    es.index(index=current_index, body=document)

forecast_table = pq.read_table(FORECAST_PATH)
df_forecast = forecast_table.to_pandas()
for _, row in df_forecast.iterrows():
    document = row.to_dict()
```

5. **Usage:** In this project, the prepared data is deployed to production systems, specifically Elasticsearch and Kibana. These tools serve as the backbone for analytics and visualization based on the data indexed. Multiple dashboards are created to cater to various analysis requirements and provide valuable insights. Users of the data, such as consumers, researchers, and statisticians can leverage the visually appealing and informative dashboards to perform a range of analyses. For example, if someone prefers to enjoy outdoor activities before or after watching a movie, they can check the weather forecast to determine the best time to visit the cinema. If it's expected to rain heavily, they might opt for an indoor movie theater instead of an outdoor cinema. On the other hand, if the weather forecast indicates clear skies and pleasant temperatures, they may choose an outdoor cinema for a unique and enjoyable experience. The versatility of Elasticsearch and Kibana empowers users to derive actionable insights from the data and make informed decisions across various domains.

## Top Rated Movies

| Movie Title ⌄ | Popularity ⌄ ↓ | Votes ⌄ |
|---|---|---|
| Reality | 123.331 | 46 |
| BlackBerry | 168.819 | 44 |
| The Tank | 91.169 | 25 |
| The Elderly | 82.359 | 17 |
| The Flash | 380.541 | 16 |
| From Black | 180.523 | 14 |
| Asteroid City | 101.347 | 3 |
| Elemental | 130.332 | 0 |
| No Hard Feelings | 91.926 | 0 |
| Indiana Jones and the Dial of Dest | 90.682 | 0 |

### Current Temperature

Paris Current Temperature

**23**

Top Rated Movies & Current Temperature

## Analysis of Top Rated Movies

(bar chart)
- Median of min_...
- Median of max...
- Median of avg_...
- Number of Rati...

Y-axis: Median of min_rating (0 – 30,000)

## Upcoming Movies

| Movies ⌄ | ↓ Release Date per da ⌄ | Number of Votes ⌄ | Popularity ⌄ |
|---|---|---|---|
| Other | 2023-06-29 | 0 | 29.201 |
| Indiana Jones and the Di | 2023-06-28 | 0 | 84.891 |
| Ruby Gillman, Teenage K | 2023-06-28 | 0 | 61.436 |
| Black Clover: Sword of th | 2023-06-16 | 0 | 45.153 |
| j-hope IN THE BOX | 2023-06-16 | 0 | 47.046 |
| No Hard Feelings | 2023-06-15 | 0 | 95.716 |
| Elemental | 2023-06-14 | 0 | 154.386 |
| Asteroid City | 2023-06-09 | 4 | 76.577 |
| Transformers: Rise of the | 2023-06-06 | 100 | 1,643.4 |

Analysis of Top Rated Movies & Upcoming Movies

(tooltip)
| Median of min_rating | 242 |
| Median of max_rating | 29,822 |
| Number of Rating | 20 |

## Weather Forecast

| Date per day ⌄ | Maximum Temperature ⌄ | Minimum Temperature ⌄ |
|---|---|---|
| 2023-06-10 | 27.95 | 17.3 |
| 2023-06-11 | 30.2 | 17.6 |
| 2023-06-12 | 29.3 | 15.9 |
| 2023-06-13 | 28.3 | 14.8 |
| 2023-06-14 | 24.7 | 12.7 |
| 2023-06-15 | 26 | 14.9 |
| 2023-06-16 | 24.6 | 13.7 |

## Trending Movies

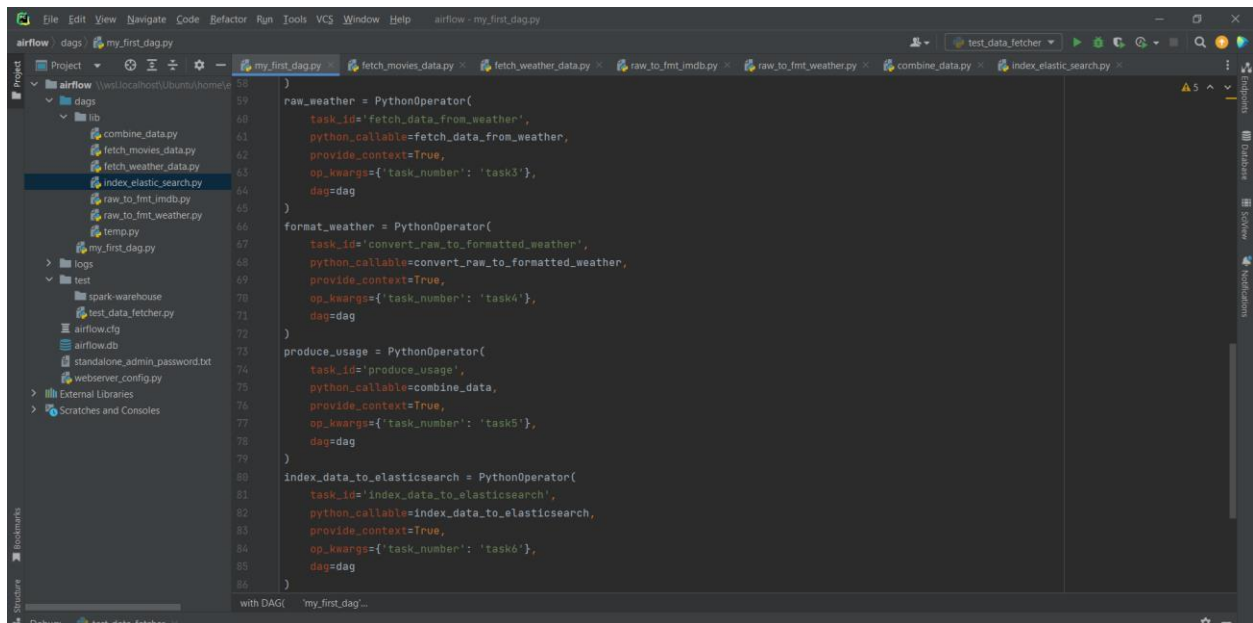| Movie Title ⌄ | Release Date per day ⌄ | Popularity ⌄ ↓ | Votes ⌄ |
|---|---|---|---|
| Spider-Man: Across the S | 2023-01-06 | 2,796.055 | 982 |
| Spider-Man: Across the S | 2023-05-31 | 2,860.755 | 925 |
| Winnie the Pooh: Blood a | 2023-01-27 | 205.719 | 631 |
| Sisu | 2023-01-27 | 954.496 | 624 |
| Marcel the Shell with Sho | 2022-06-24 | 168.785 | 284 |
| The Wandering Earth II | 2023-01-22 | 451.908 | 261 |
| Other | 2023-01-05 | 134.974 | 170 |
| Other | 2023-01-13 | 69.139 | 118 |
| Transformers: Rise of the | 2023-06-06 | 1,643.4 | 100 |

Weather Forecast & Trending Movies

## Directed Acyclic Graphs - DAGs

We use Apache Airflow to create DAGs (Directed Acyclic Graphs) to manage workflow orchestration of our data. The tasks and dependencies are defined in python (Figure 10 and 11) and Airflow manages the scheduling and execution. The DAG that reflects our pipeline architecture is shown in Figure 12.



Figure 10: DAG Phyton Code i
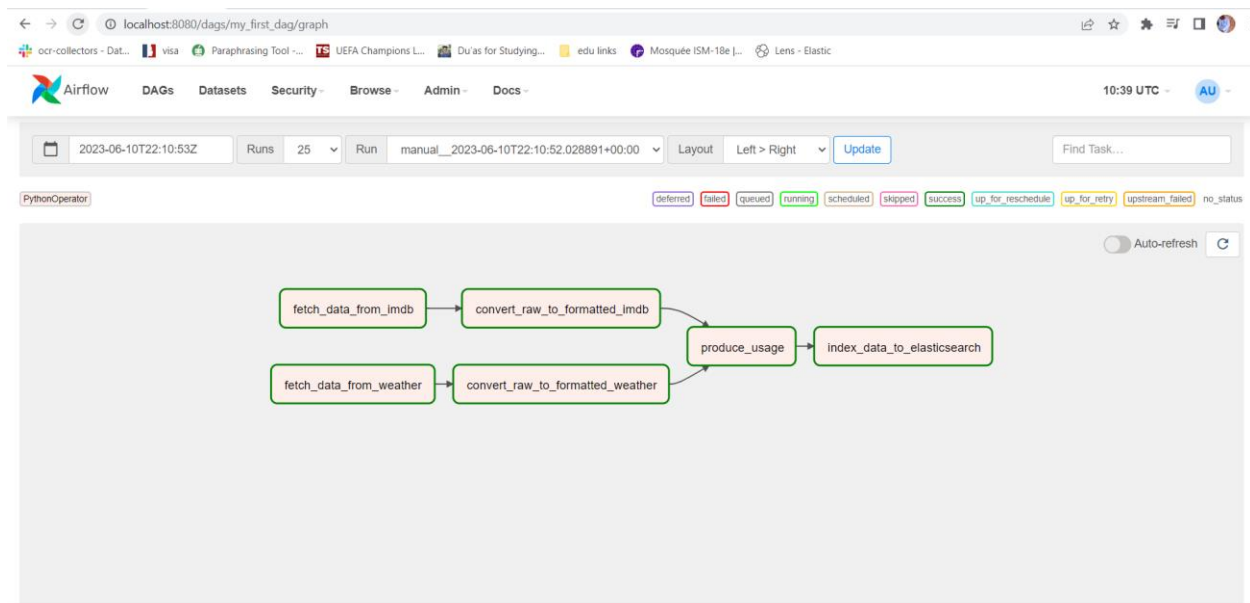


Figure 11: DAG Phyton Code ii

Figure 12 : DAG Airflow Graph View

**Conclusion**

Throughout the implementation of this project, we have gained extensive knowledge and valuable experience in various technologies. Despite facing challenges along the way, the outcome has been highly successful, providing us with in-depth learnings and honing our skills.

The dynamic and comprehensive visualization of cinema and weather data offers significant benefits to consumers, researchers, statisticians, data scientists, and others, empowering them to make well-informed decisions, perform comparisons, and apply filters based on their specific needs. By considering trending movies, upcoming releases, top-rated films, user ratings, and weather forecasts for the upcoming week, stakeholders can leverage these insights for powerful recommendations.

This project has allowed us to broaden our understanding of key concepts such as big data, data lakes, data transformation, visualization, streaming, data analytics, and workflow. We have become proficient in utilizing various tools and technologies, including Apache Airflow, Apache Spark, cloud-based Elasticsearch and Kibana, Python, Panda, Parquet, and JSON, among others.