# Big Data Project

**TEAM MEMBERS**

**ASMAU AMINU 62709**

**YUVASRI BALU 62780**

## Table of Content

## 1.  INTRODUCTION

The Big Data Project is an ambitious undertaking that aims to demonstrate the utilization of Apache Airflow, a powerful workflow management platform, in conjunction with other Big Data technologies to create an end-to-end data pipeline. This project involves extracting data from external sources, preparing and formatting the data, performing data analysis using Apache Spark, and ultimately generating valuable insights.

The project will be divided into several episodes, each focusing on specific aspects of the overall data pipeline. The episodes will cover the fundamentals of Airflow, hands-on exercises to set up and familiarize oneself with Airflow, creating custom Directed Acyclic Graphs (DAGs) to orchestrate tasks, and integrating the project with Git for version control.

In later episodes, the project will delve into extracting data from mediaStack and exchangerate APIs, preparing and formatting the data into structured layers, and combining multiple datasets for comprehensive analysis. Apache Spark, a distributed computing system, will be utilized to process and analyze the data, generating meaningful insights and output.

hroughout the report, we will provide detailed documentation, step-by-step instructions, and code samples to guide you through each episode of the project.

## Data Source and API

Free foreign exchange, crypto rates &EU VAT Rates API

https://exchangerate.host/#/#docs

It shows the exchange rate of many currencies with euro as a base.

News API

https://mediastack.com/documentation

It contains daily news with author, tittle, description, date.

Structure of the Datalake

`datalake/`: The root directory of the data lake.
`raw/`: Contains the raw data obtained directly from the APIs.
`exchangerates/`: Stores the raw data from the Exchangerates API.
`money/`: Stores the exchange rate data obtained from the Exchangerates API.
`<current_day>/`: Subdirectory for each day the data is fetched.
`<raw_data_files>`: Raw data files in JSON format (e.g., `money.json`).
`mediastack/`: Stores the raw data from the Mediastack API.
`news/`: Stores the news data obtained from the Mediastack API.
`<current_day>/`: Subdirectory for each day the data is fetched.
`<raw_data_files>`: Raw data files in JSON format (e.g., `news.json`).
`formatted/`: Contains the processed and formatted data.
`exchangerates/`: Stores the formatted data of exchange rates.
`money/`: Stores the formatted exchange rate data.
`<current_day>/`: Subdirectory for each day the data is processed.
`<formatted_data_files>`: Formatted data files in Parquet format (e.g., `money.snappy.parquet`).
`mediastack/`: Stores the formatted data of news.

**news/**: Stores the formatted news data.
**<current_day>/**: Subdirectory for each day the data is processed.
**<formatted_data_files>**: Formatted data files in Parquet format (e.g., **news.snappy.parquet**).
**usage/**: Contains the analyzed and indexed data for usage analysis.
**exchangeAnalysis/**: Stores the analyzed data related to exchange rates.
**ExchangedStatistics/**: Stores the analyzed exchange rate statistics.
**<current_day>/**: Subdirectory for each day the data is analyzed.
**<analyzed_data_files>**: Analyzed data files in Parquet format (e.g., **res.snappy.parquet**).
**newsAnalysis/**: Stores the analyzed data related to news.
**NewsAnalysis/**: Stores the analyzed news data.
**<current_day>/**: Subdirectory for each day the data is analyzed.
**<analyzed_data_files>**: Analyzed data files in Parquet format (e.g., **res.snappy.parquet**).

## Pipeline

- The pipeline of this project involves several steps to collect, process, and analyze data. Let's walk through the pipeline step by step:
- **Fetching Data from Exchangerates API:**
- The fetch_data_exchangerates function is executed as a task in the pipeline.
- This function makes an HTTP request to the Exchangerates API to retrieve the latest exchange rate data.
- The data is downloaded and stored in the raw data folder of the data lake.
- **Converting Raw Exchangerates Data to Formatted Data:**
- The convert_raw_to_formatted_exchangerates function is executed as a task in the pipeline.
- It takes the raw data obtained from the Exchangerates API and processes it to create a formatted representation.
- The raw data in JSON format is converted to a Pandas DataFrame and then saved as a Parquet file in the formatted data folder of the data lake.
- **Fetching Data from Mediastack API:**

- The fetch_data_from_mediastack function is executed as a task in the pipeline.
- It makes an HTTP request to the Mediastack API to retrieve news articles and headlines.
- The data is downloaded and stored in the raw data folder of the data lake.
- **Converting Raw Mediastack Data to Formatted Data:**
- The convert_raw_to_formatted_mediatstack function is executed as a task in the pipeline.
- It processes the raw news data obtained from the Mediastack API and converts it to a formatted representation.
- The raw JSON data is transformed into a Pandas DataFrame and then saved as a Parquet file in the formatted data folder of the data lake.
- **Combining and Analyzing Data:**
- The combine_data_api function is executed as a task in the pipeline.
- It combines the formatted exchange rate data and news data from the data lake to perform data analysis.
- The Spark framework is used to analyze the data, extract relevant information, and generate insights.

**Indexing Data into Elasticsearch:**

- The index_elasticsearch function is executed as a task in the pipeline.
- It connects to Elasticsearch and indexes the analyzed data into appropriate indices.
- The data is transformed into JSON records and bulk indexed into Elasticsearch for efficient searching and querying.

## 2. Create our first DAGS

## 3. Big Data Project - Extract source from Media Stack



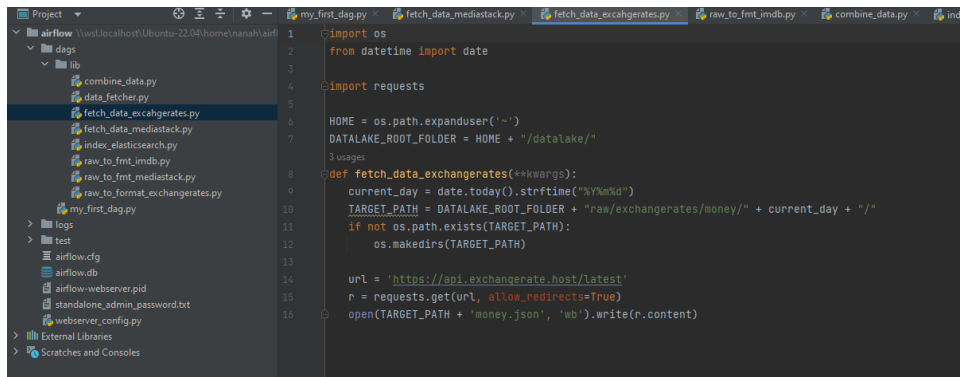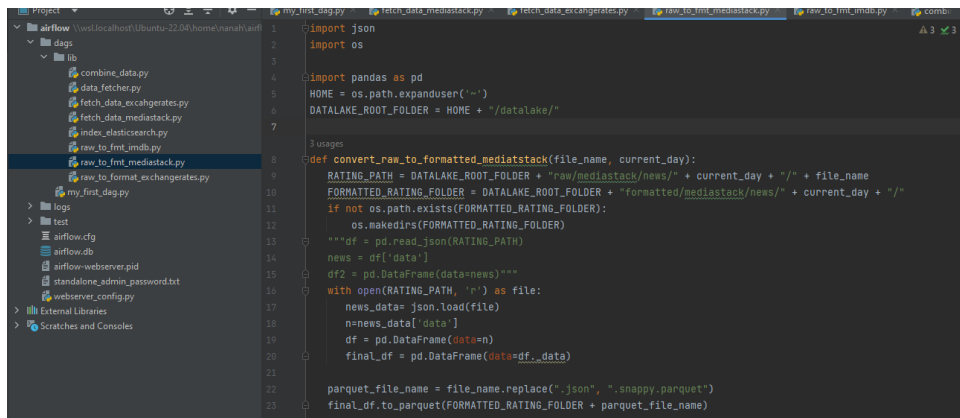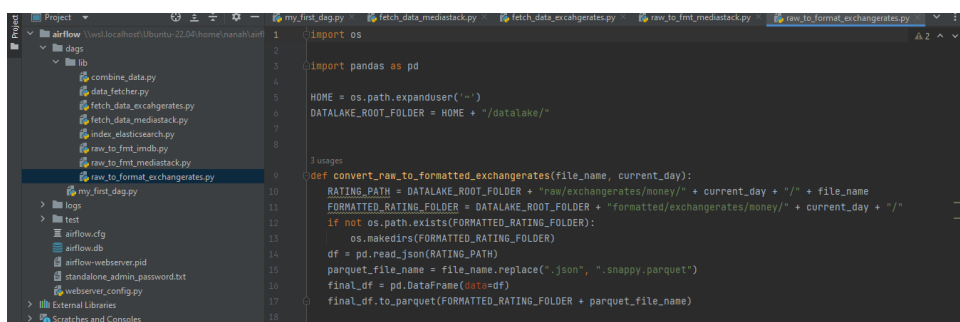## 4. Big Data Project - Extract source from Exchange rate

```python
import os

from datetime import date

import requests


HOME = os.path.expanduser('~')
DATALAKE_ROOT_FOLDER = HOME + "/datalake/"

3 usages
def fetch_data_exchangerates(**kwargs):
    current_day = date.today().strftime("%Y%m%d")
    TARGET_PATH = DATALAKE_ROOT_FOLDER + "raw/exchangerates/money/" + current_day + "/"
    if not os.path.exists(TARGET_PATH):
        os.makedirs(TARGET_PATH)


    url = 'https://api.exchangerate.host/latest'
    r = requests.get(url, allow_redirects=True)
    open(TARGET_PATH + 'money.json', 'wb').write(r.content)
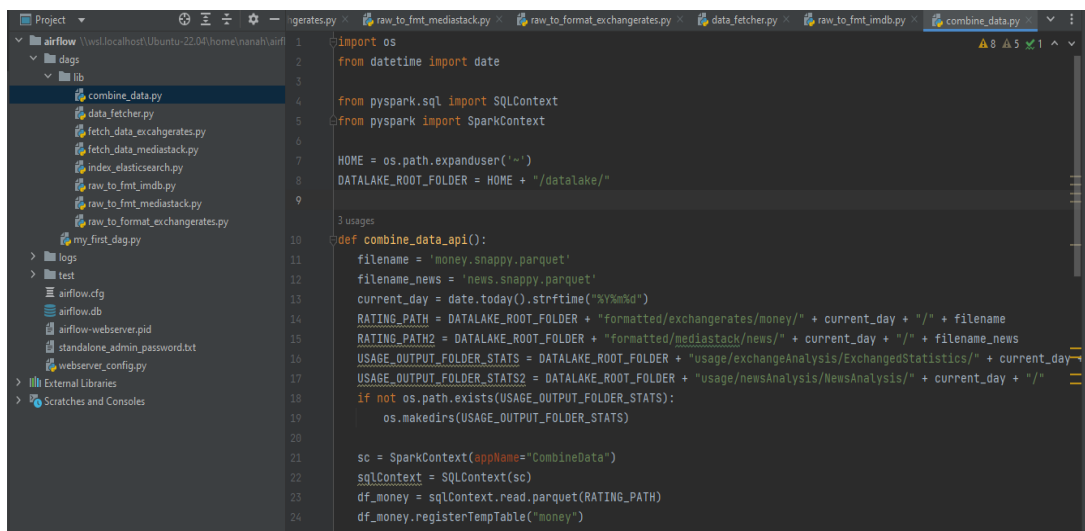```

## 5. Big Data Project - Prepare formatted data



```python
import json
import os

import pandas as pd
HOME = os.path.expanduser('~')
DATALAKE_ROOT_FOLDER = HOME + "/datalake/"


3 usages
def convert_raw_to_formatted_mediatack(file_name, current_day):
    RATING_PATH = DATALAKE_ROOT_FOLDER + "raw/mediastack/news/" + current_day + "/" + file_name
    FORMATTED_RATING_FOLDER = DATALAKE_ROOT_FOLDER + "formatted/mediastack/news/" + current_day + "/"
    if not os.path.exists(FORMATTED_RATING_FOLDER):
        os.makedirs(FORMATTED_RATING_FOLDER)
    """df = pd.read_json(RATING_PATH)
    news = df['data']
    df2 = pd.DataFrame(data=news)"""
    with open(RATING_PATH, 'r') as file:
        news_data= json.load(file)
        n=news_data['data']
        df = pd.DataFrame(data=n)
        final_df = pd.DataFrame(data=df._data)

    parquet_file_name = file_name.replace(".json", ".snappy.parquet")
    final_df.to_parquet(FORMATTED_RATING_FOLDER + parquet_file_name)
```



```python
import os

import pandas as pd

HOME = os.path.expanduser('~')
DATALAKE_ROOT_FOLDER = HOME + "/datalake/"


3 usages
def convert_raw_to_formatted_exchangerates(file_name, current_day):
    RATING_PATH = DATALAKE_ROOT_FOLDER + "raw/exchangerates/money/" + current_day + "/" + file_name
    FORMATTED_RATING_FOLDER = DATALAKE_ROOT_FOLDER + "formatted/exchangerates/money/" + current_day + "/"
    if not os.path.exists(FORMATTED_RATING_FOLDER):
        os.makedirs(FORMATTED_RATING_FOLDER)
    df = pd.read_json(RATING_PATH)
    parquet_file_name = file_name.replace(".json", ".snappy.parquet")
    final_df = pd.DataFrame(data=df)
    final_df.to_parquet(FORMATTED_RATING_FOLDER + parquet_file_name)
```

## 6. Big Data Project - Prepare combined data



```python
import os

from datetime import date

from pyspark.sql import SQLContext
from pyspark import SparkContext


HOME = os.path.expanduser('~')
DATALAKE_ROOT_FOLDER = HOME + "/datalake/"


3 usages
def combine_data_api():
    filename = 'money.snappy.parquet'
    filename_news = 'news.snappy.parquet'
    current_day = date.today().strftime("%Y%m%d")
    RATING_PATH = DATALAKE_ROOT_FOLDER + "formatted/exchangerates/money/" + current_day + "/" + filename
    RATING_PATH2 = DATALAKE_ROOT_FOLDER + "formatted/mediastack/news/" + current_day + "/" + filename_news
    USAGE_OUTPUT_FOLDER_STATS = DATALAKE_ROOT_FOLDER + "usage/exchangeAnalysis/ExchangedStatistics/" + current_day
    USAGE_OUTPUT_FOLDER_STATS2 = DATALAKE_ROOT_FOLDER + "usage/newsAnalysis/NewsAnalysis/" + current_day + "/"
    if not os.path.exists(USAGE_OUTPUT_FOLDER_STATS):
        os.makedirs(USAGE_OUTPUT_FOLDER_STATS)


    sc = SparkContext(appName="CombineData")
    sqlContext = SQLContext(sc)
    df_money = sqlContext.read.parquet(RATING_PATH)
    df_money.registerTempTable("money")
```

# 7. Big Data Project - Ingest into EK

```python
#Read the Parquet file into a DataFrame
table = pq.read_table(RATING_PATH)
table2 = pq.read_table(RATING_PATH2)
df = table.to_pandas()
df2 = table2.to_pandas()

# Convert DataFrame to JSON records
json_data = df.to_json(orient='records')
json_data2 = df2.to_json(orient='records')

# Parse JSON records into a list of dictionaries
documents = json.loads(json_data)
documents2 = json.loads(json_data2)

# Bulk index the documents into Elasticsearch
response = bulk(es, documents, index=index_name)
response2 = bulk(es, documents2, index=index_name2)

# Check if the indexing operation was successful
if response[0] > 0:
    print("Indexing successful!")
else:
    print("Indexing failed.")

# Retrieve information about the index
response = es.indices.get(index=index_name)
response2 = es.indices.get(index=index_name2)

# Print the retrieved information
print(response)
print(response2)
```
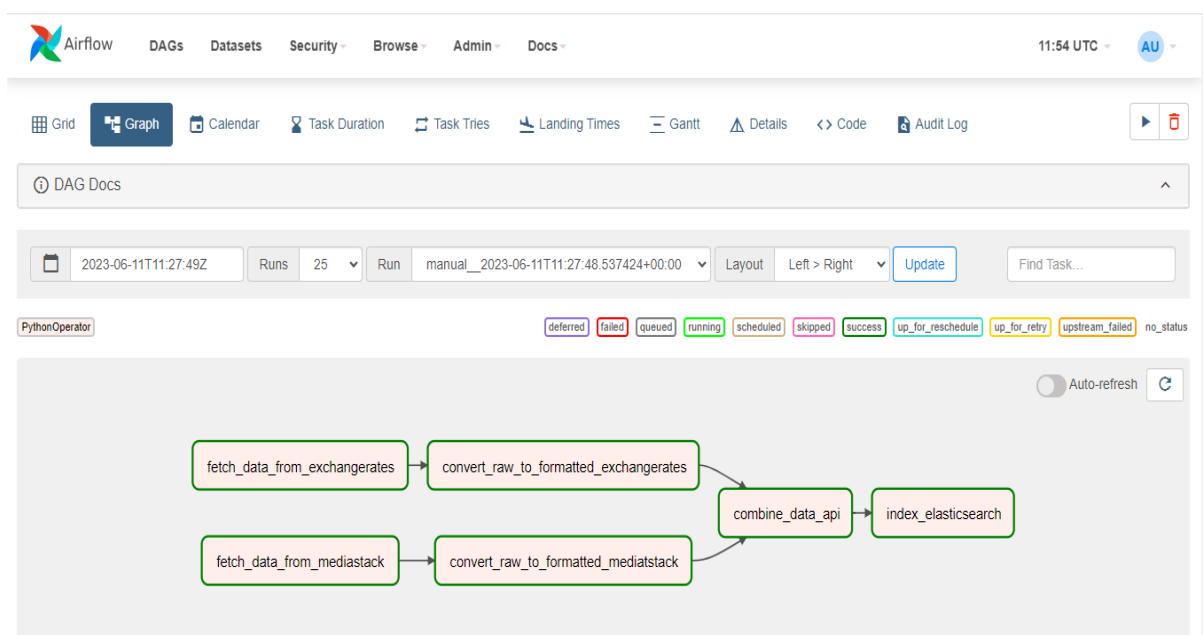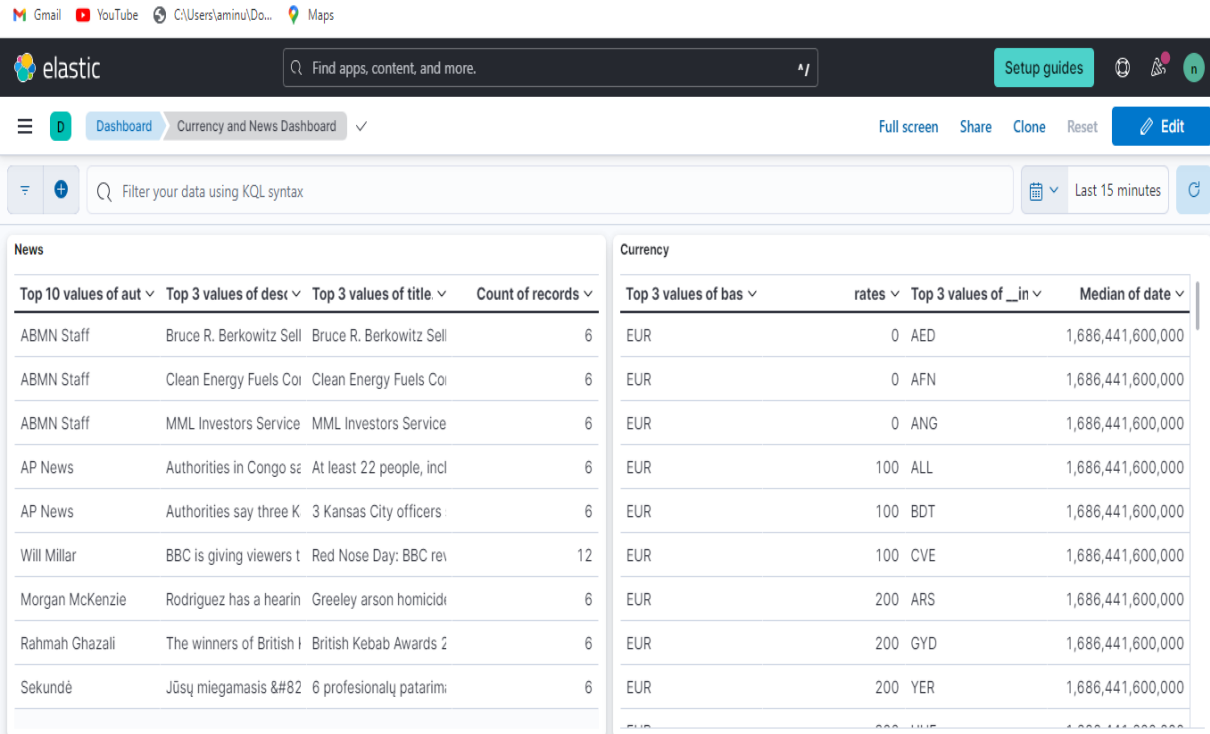
## VIsualization

DASHBOARD



# Conclusion

In conclusion, this project demonstrates a data pipeline for collecting, processing, analyzing, and indexing data from two different sources: the Exchangerates API and the Mediastack API. The pipeline is designed to run periodically, fetching the latest data from the APIs and performing various data processing tasks.