# Princess Sumaya University for Technology

# Data Engineering Course

# Assignment 1

# Dr Ibrahim Abu Alhaol

By Waed Alsawarieh , 20208020

# Question 1

Q1: Provide similar to AirFlow implementation but using NiFi and provide the GitHub repo with all dependencies and detailed REAME.MD and PPT presentation on how to run your workflow.

# 1.1 Get Data CSV From Input Directory in NIFI Container

# 1.2 UpdateAttribute Processor to update attribute from data.csv to data.json

# 1.3 ConvertRecord Processor

Converts records from one data format to another using configured Record Reader and Record Write Controller Services.

# 1.4 PutFile Processor

Writes the contents of a FlowFile to the local file system



**Processor Details**

▶ Running                                                    ⚙ STOP & CONFIGURE

| SETTINGS | SCHEDULING | PROPERTIES | COMMENTS |

**Required field**

| Property | | Value |
|---|---|---|
| **Directory** | ❓ | /opt/nifi/output |
| **Conflict Resolution Strategy** | ❓ | replace |
| **Create Missing Directories** | ❓ | true |
| Maximum File Count | ❓ | No value set |
| Last Modified Time | ❓ | No value set |
| Permissions | ❓ | No value set |
| Owner | ❓ | No value set |
| Group | ❓ | No value set |

OK

```
waedas@waedas-Inspiron-5584:~/Desktop/DEAssignment1/nifi$ docker exec -it a87ffd8b8020 bash
nifi@a87ffd8b8020:/opt/nifi/nifi-current$ cd ../output/
nifi@a87ffd8b8020:/opt/nifi/output$ ls
data.json
nifi@a87ffd8b8020:/opt/nifi/output$ █
```

Validate data.json

# Question 2

Q2: Provide Similar to Airflow implementation but with csv file is extracted from Postgresql table and the produced json file is pushed to MongoDB database. Provide Github repo with all dependencies and detailed REAME.MD and PPT presentation how to run your workflow.

# Docker Compose File- Services

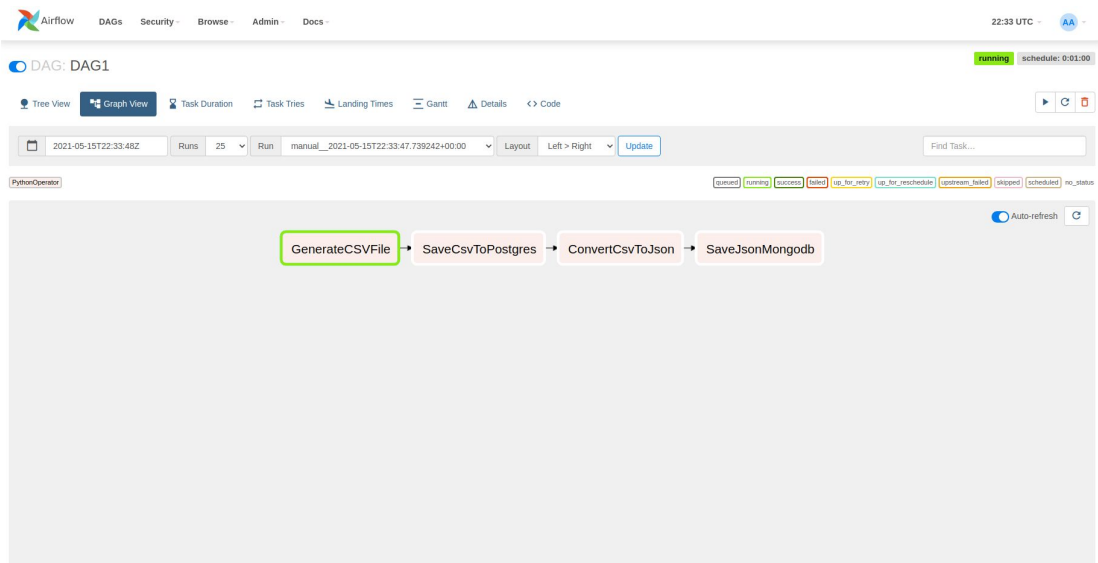1.**Apache Airflow**

2.**Postgresql**

3.**pgAdmin**

4.**mongo**-**express**: web-based **MongoDB** admin interface.

5.**mongo** : MongoDB document database.

# 2.1 Generate CSV file from faker

```python
def GenerateCSV():
    output = open('data.csv', 'w')
    fake = Faker()
    header = ['name', 'age', 'street', 'city', 'state', 'zip', 'lng', 'lat']
    mywriter = csv.writer(output)
    mywriter.writerow(header)
    for r in range(10):
        row = [fake.name(), fake.random_int(min=18, max=80, step=1),
               fake.street_address(), fake.city(), fake.state(),
               fake.zipcode(), fake.longitude(), fake.latitude()]
        print(row)
        mywriter.writerow(row)

    output.close()
    DF = pd.read_csv('data.csv')
    DF.to_csv(AIRFLOW_HOME + '/dags/dataframe.csv', index=False)
```
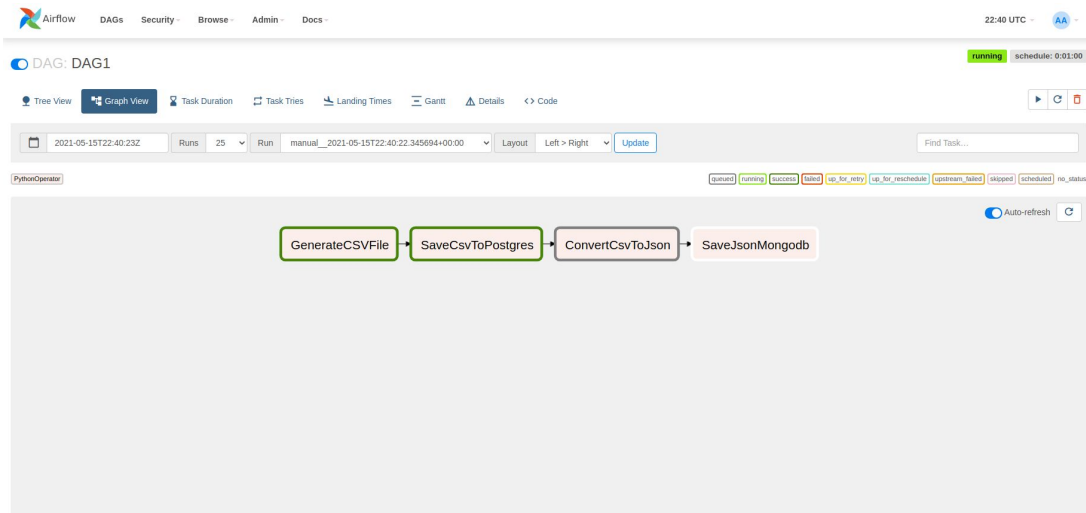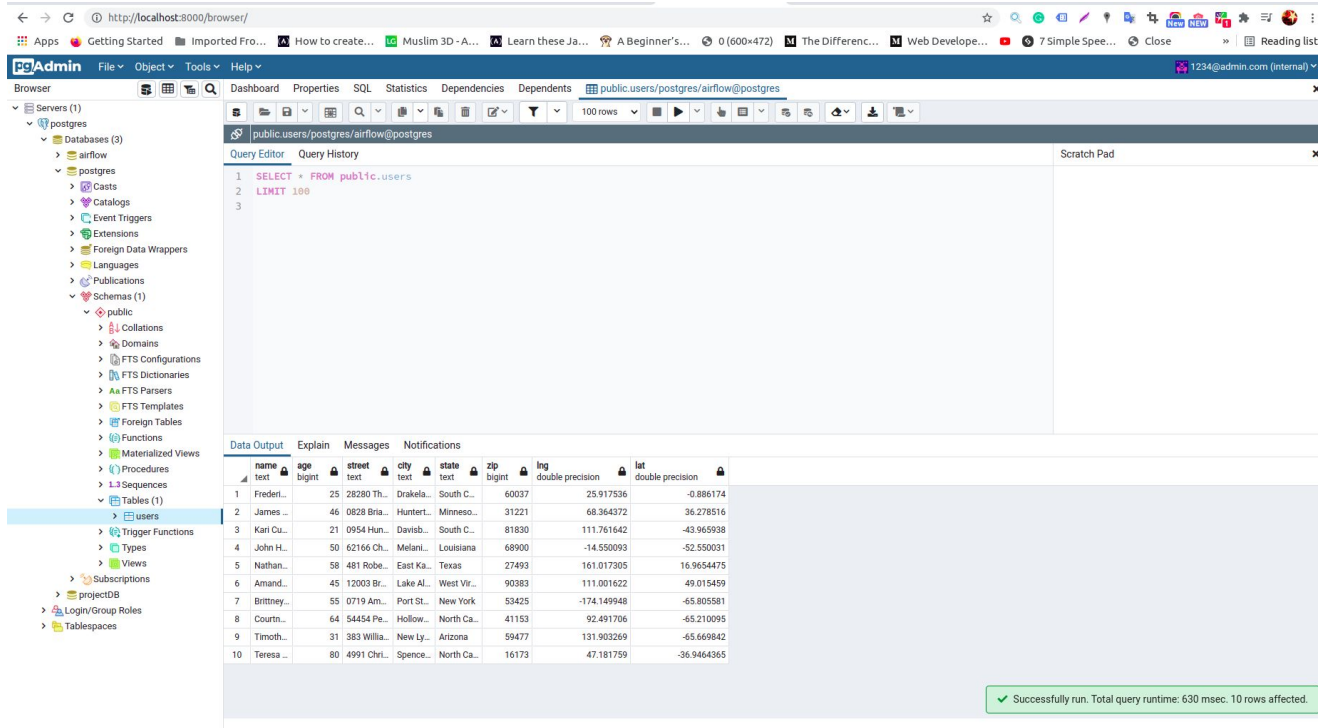
# 2.2 Save CSV in Postgresql Database

```python
# config varabiles
host = Variable.set("host", "postgres")
user = Variable.set("user", "airflow")
password = Variable.set("password", "airflow")
port = Variable.set("port", '5432')
database = Variable.set("database", 'postgres')
AIRFLOW_HOME = os.getenv('AIRFLOW_HOME')

def SaveCsvToPostgres():
    host = Variable.get('host')
    user = Variable.get('user')
    password = Variable.get('password')
    port = Variable.get('port')
    database = Variable.get('database')
    engine = create_engine(
        f'postgresql://{user}:{password}@{host}:{port}/{database}')
    print("Airflow Database Tables :- ", engine.table_names())
    DF = pd.read_csv(AIRFLOW_HOME + '/dags/dataframe.csv')
    # push table
    DF.to_sql('users', engine, if_exists='replace', index=False)
```
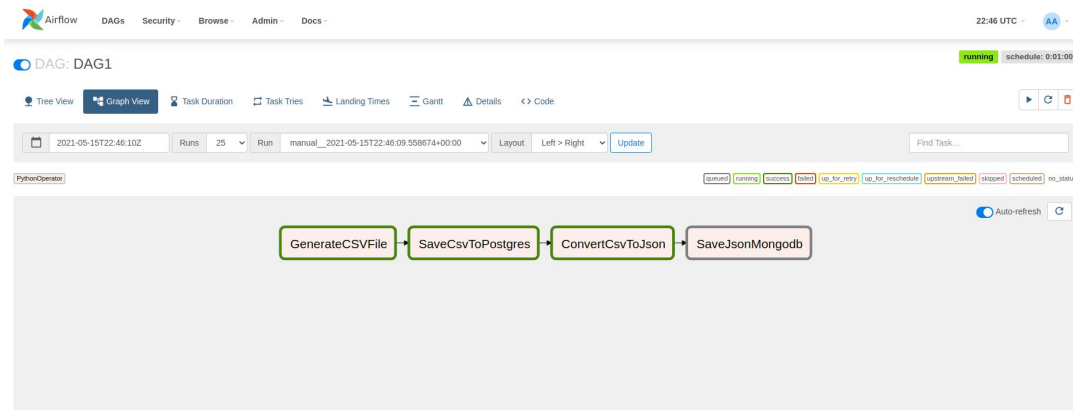
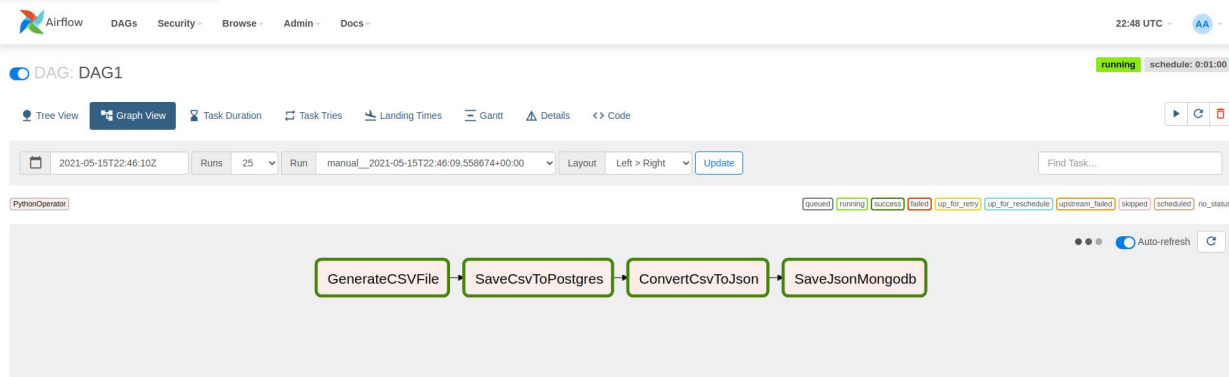# 2.3 Save CSV in Postgresql Database

# 2.4 Convert To JSON

```python
def ConvertCsvToJson():
    # read from postgres
    host = Variable.get('host')
    user = Variable.get('user')
    password = Variable.get('password')
    port = Variable.get('port')
    database = Variable.get('database')
    engine = create_engine(
        f'postgresql://{user}:{password}@{host}:{port}/{database}')
    DF2 = pd.read_sql("SELECT * FROM users", engine)

    for i, r in DF2.iterrows():
        print(r['name'])

    DF2.to_json(AIRFLOW_HOME + '/dags/fromAirflow.json', orient='records')
```

# 2.4 Save JSON file in MongoDB

```python
def SaveJsonMongodb():
    from pymongo import MongoClient
    client = MongoClient('mongo:27017',
                         username='root',
                         password='example')
    db = client['users']
    # Create Collection
    usersInfo = db.usersInfo
    with open(AIRFLOW_HOME + '/dags/fromAirflow.json') as f:
        users = json.load(f)
    # Push documents to collection
    for key in users:
        usersInfo.insert_one(key)
```

# 2.4 Save Json file in MongoDB

# Research

Explain all of the research you've done about this issue/challenge.

What was the goal of your research? Be sure to explain how you found it and anyone who might have helped you!