# ASSIGNMENT – 3 (Apache Airflow)

# TRAINEE NAME: Swathi Baskaran
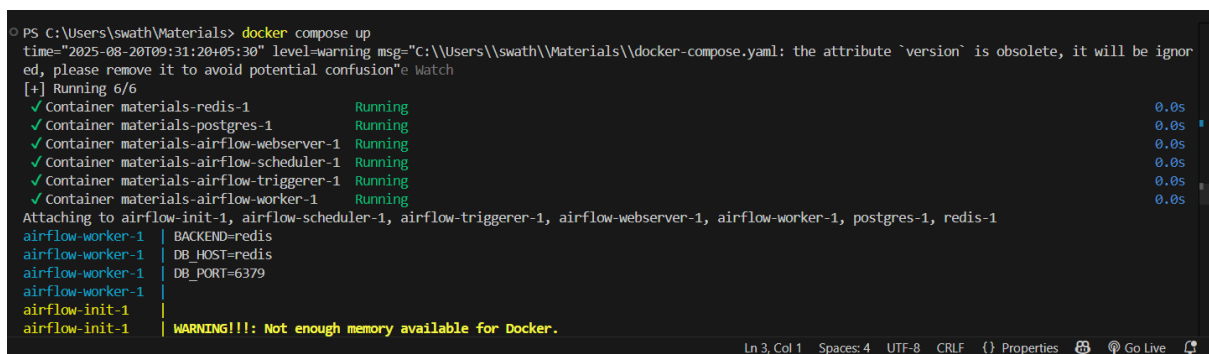
## Setting up airflow-init in VS Code



## Starting up all services

# Setting up a new connection



# Entering the connection details



# Connection successfully created

# Create tables for staging and final data

```python
from airflow.providers.postgres.operators.postgres import PostgresOperator

create_employees_table = PostgresOperator(
    task_id="create_employees_table",
    postgres_conn_id="tutorial_pg_conn",
    sql="""
        CREATE TABLE IF NOT EXISTS employees (
            "Serial Number" NUMERIC PRIMARY KEY,
            "Company Name" TEXT,
            "Employee Markme" TEXT,
            "Description" TEXT,
            "Leave" INTEGER
        );""",
)

create_employees_temp_table = PostgresOperator(
    task_id="create_employees_temp_table",
    postgres_conn_id="tutorial_pg_conn",
    sql="""
        DROP TABLE IF EXISTS employees_temp;
        CREATE TABLE employees_temp (
            "Serial Number" NUMERIC PRIMARY KEY,
            "Company Name" TEXT,
            "Employee Markme" TEXT,
            "Description" TEXT,
            "Leave" INTEGER
        );""",
)
```

# Load data into the staging table

```python
import os
import requests
from airflow.decorators import task
from airflow.providers.postgres.hooks.postgres import PostgresHook


@task
def get_data():
    # NOTE: configure this as appropriate for your airflow environment
    data_path = "/opt/airflow/dags/files/employees.csv"
    os.makedirs(os.path.dirname(data_path), exist_ok=True)

    url = "https://raw.githubusercontent.com/apache/airflow/main/airflow-core/docs/tutorial/pipeline_example.csv"

    response = requests.request("GET", url)

    with open(data_path, "w") as file:
        file.write(response.text)

    postgres_hook = PostgresHook(postgres_conn_id="tutorial_pg_conn")
    conn = postgres_hook.get_conn()
    cur = conn.cursor()
    with open(data_path, "r") as file:
        cur.copy_expert(
            "COPY employees_temp FROM STDIN WITH CSV HEADER DELIMITER AS ',' QUOTE '\"'",
            file,
        )
    conn.commit()
```

## Merge and clean the data

```python
from airflow.decorators import task
from airflow.providers.postgres.hooks.postgres import PostgresHook


@task
def merge_data():
    query = """
        INSERT INTO employees
        SELECT *
        FROM (
            SELECT DISTINCT *
            FROM employees_temp
        ) t
        ON CONFLICT ("Serial Number") DO UPDATE
        SET
            "Employee Markme" = excluded."Employee Markme",
            "Description" = excluded."Description",
            "Leave" = excluded."Leave";
    """
    try:
        postgres_hook = PostgresHook(postgres_conn_id="tutorial_pg_conn")
        conn = postgres_hook.get_conn()
        cur = conn.cursor()
        cur.execute(query)
        conn.commit()
        return 0
    except Exception as e:
        return 1
```

## Defining the DAG

```python
import datetime
import pendulum
import os
import requests

from airflow.decorators import dag, task
from airflow.providers.postgres.hooks.postgres import PostgresHook
from airflow.providers.postgres.operators.postgres import PostgresOperator




@dag(
    dag_id="process_employees",
    schedule="0 0 * * *",  # daily at midnight
    start_date=pendulum.datetime(2021, 1, 1, tz="UTC"),
    catchup=False,
    dagrun_timeout=datetime.timedelta(minutes=60),
)
def ProcessEmployees():
    # Step 1: Create main employees table
    create_employees_table = PostgresOperator(
        task_id="create_employees_table",
        postgres_conn_id="tutorial_pg_conn",
        sql="""
            CREATE TABLE IF NOT EXISTS employees (
                "Serial Number" NUMERIC PRIMARY KEY,
                "Company Name" TEXT,
                "Employee Markme" TEXT,
                "Description" TEXT,
                "Leave" INTEGER
            );""",
    )

    # Step 2: Create staging table
    create_employees_temp_table = PostgresOperator(
        task_id="create_employees_temp_table",
        postgres_conn_id="tutorial_pg_conn",
```

# Finding the DAG in the Apache UI

◯ DAG: process_employees                                    Schedule: 0 0 * * * ⓘ   Next Run: 2025-08-20, 00:00

▦ Grid    ▫ Graph    📅 Calendar    ⧗ Task Duration    ⇄ Task Tries    ⬎ Landing Times    ☰ Gantt    ⚠ Details    <> Code    🔒 Audit Log                ▶

20-08-2025 04:15:51 AM 📅    25 ⌄    All Run Types ⌄    All Run States ⌄    **Clear Filters**

deferred | failed | queued | removed | restarting | running | scheduled | shutdown | skipped | success | up_for_reschedule | up_for_retry | upstream_failed | no_sta

Auto-refresh ◯                                                    ⇥☰

|  | DAG |  |
|---|---|---|
|  | **process_employees** |  |

Duration

00:03:19

00:01:39

00:00:00

create_employees_table
create_employees_temp_table
get_data
merge_data

**DAG Details**

| DAG Runs Summary | |
|---|---|
| Total Runs Displayed | 2 |
| ■ Total success | 2 |
| First Run Start | 2025-08-19, 06:23:52 UTC |
| Last Run Start | 2025-08-20, 03:59:29 UTC |
| Max Run Duration | 00:03:19 |
| Mean Run Duration | 00:01:44 |

---

◯ DAG: process_employees                          success   Schedule: 0 0 * * * ⓘ   Next Run: 2025-08-19, 00:0

▦ Grid    ▫ Graph    📅 Calendar    ⧗ Task Duration    ⇄ Task Tries    ⬎ Landing Times    ☰ Gantt    ⚠ Details    <> Code    🔒 Audit Log                ▶

📅 2025-08-18T00:00:01Z    Runs 25 ⌄    Run scheduled__2025-08-18T00:00:00+00:00 ⌄    Layout Left > Right ⌄    **Update**          Find Task…

PostgresOperator | _PythonDecoratedOperator          deferred | failed | queued | removed | restarting | running | scheduled | shutdown | skipped | success | up_for_reschedule | up_for_retry | upstream_failed | no_

Auto-refresh ◯