

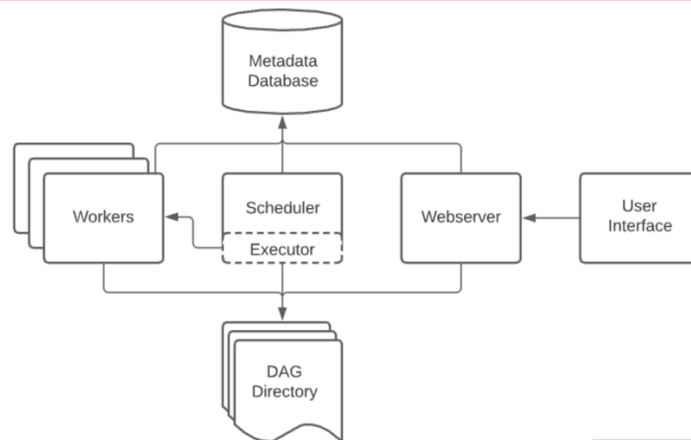
HƯỚNG DẪN APACHE AIRFLOW

THEORETICAL BASIS

1. Definition

Apache airflow là một nền tảng giúp hỗ trợ việc lập lịch và quản lý các workflow thông qua việc lập trình

2. Overview Architecture



I. DOWNLOAD AND INSTALL

Method 1: Direct Installation on Local Machine

a) Installation

- Install Airflow via Pip:

```
pip install apache-airflow
```

- Initialize database

```
airflow db init
```

- Create user

```
airflow users create --username [tên đăng nhập] --password [mật khẩu] --firstname [Tên] --lastname [Họ] --role Admin --email [email]
```

- Configuration: customize the settings in the 'airflow.cfg' file

b) Run

- Initialize web server

```
airflow webserver -p 8080
```

- Initialize Scheduler (in new terminal)

```
airflow scheduler
```

Method 2: Use Docker

Giảng viên biên soạn: Phạm Thị Xuân Hiền

- Open new terminal

```
docker run -it --name demo1 -v
/Users/macintoshhd/Desktop/apache_airflow/dags:/opt/airflow/dags -p 8080:8080 -d
apache/airflow:2.8.4-python3.9 bash
```

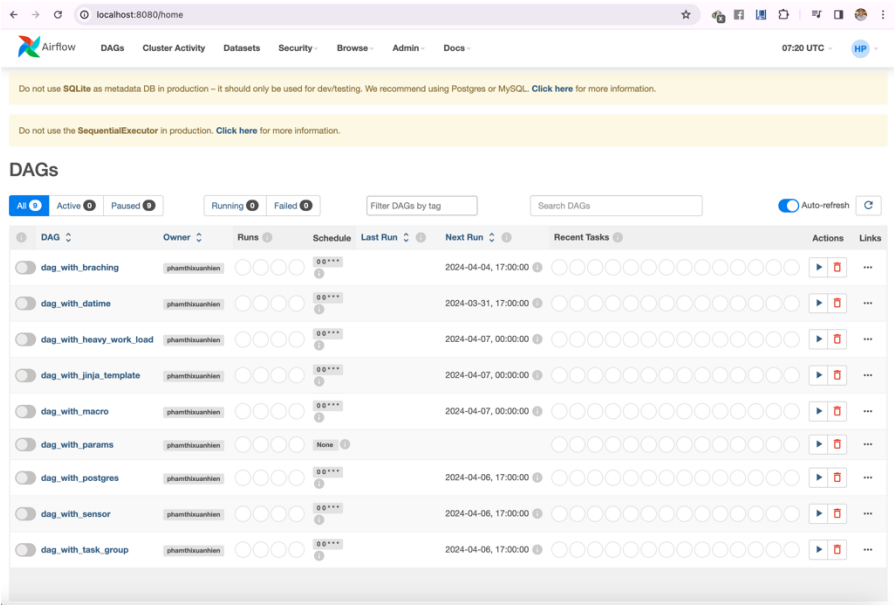
- Open Docker Desktop

<input type="checkbox"/>	Name	Image	Status	CPU (%)	Port(s)	Last started	Actions
<input type="checkbox"/>	<div>demo1</div> <div>4a12a71236fd</div>	apache/airflow:2.8.4-python3.9	Exited (255)	N/A	8080:8080	2 days ago	<div><div></div><div></div><div></div></div>

- Copy all files from the dags local directory to the dags folder on Apache Airflow

```
bash -c '(airflow db init && airflow users create --username admin --password admin
--firstname Hien --lastname Pham --role Admin --email phamthixuanhien@iuh.edu.vn);
airflow webserver & airflow scheduler'
```

- Open “localhost:8080” in Google Chrome



link ref: <https://airflow.apache.org/docs/apache-airflow/1.10.4/start.html>

II. INTRODUCE

1. DAG with Branch

```
dags > dag_with_branching.py > ...
1 import datetime as dt
2 import pendulum
3 from airflow import DAG
4 from airflow.decorators import task, task_group
5 from airflow.operators.empty import EmptyOperator
6 from airflow.utils.trigger_rule import TriggerRule
7
8 # utc = pendulum.timezone("UTC")
9 hcm_tz = pendulum.timezone("Asia/Ho_Chi_Minh")
10
11 default_args = {
12     "owner": "phamthixuanhien",
13     # "start_date": dt.datetime(2024, 4, 5, tzinfo=hcm_tz),
14     "start_date": dt.datetime(year=2024, month=4, day=5, tzinfo=hcm_tz),
15     # "start_date": dt.datetime(year=2024, month=4, day=1, tzinfo=utc),
16     "retries": 1,
17     "retry_delay": dt.timedelta(minutes=5),
18     "depend_on_past": True,
19     "email": ["phamthixuanhien@iuh.edu.vn"]
20 }
21
22 @task.branch(task_id="branching")
23 def task_branching():
24     return "task_1"
25
26 with DAG(
27     "dag_with_branching",
28     schedule_interval="0 0 * * *",
29     default_args=default_args,
30     catchup=True,
31 ) as dag:
32     start = EmptyOperator(task_id="start")
33
34     branch = task_branching()
35
36     task_1 = EmptyOperator(task_id="task_1")
37     task_2 = EmptyOperator(task_id="task_2")
38
39     end = EmptyOperator(task_id="end",
40                         trigger_rule=TriggerRule.NONE_FAILED_MIN_ONE_SUCCESS
41                         )
42
43     # Flow
44     start >> branch >> [ task_1 , task_2 ] >> end
45
46
```

with DAG(...): Khởi tạo 1 DAG với tên là **dag_with_branching**
Schedule_interval= "0 0 * * *": DAG này chạy hàng ngày vào lúc nửa đêm

Retries:1 - số lần thử lại là 1

Retry_delay: khoảng thời gian chờ giữa các lần thử lại: 5 phút

start = EmptyOperator(task_id="start"): Tạo một task trống với ID là 'start'

branch = task_branching(): Một task cho phép thực hiện nhánh điều kiện (chưa rõ hành động cụ thể vì nó phụ thuộc vào định nghĩa của task_branching())

trigger_rule=TriggerRule.NONE_FAILED_MIN_ONE_SUCCESS: nghĩa là task này chỉ thực thi khi ít nhất một trong các tasks trước nó thành công và không có task nào thất bại.

Định nghĩa luồng công việc (workflow): **start** chạy đầu tiên >> tiếp theo là **branch**, sau đó dựa trên kết quả của **branch**, một trong hai (hoặc cả hai) **task_1**, **task_2** sẽ được chạy >> cuối cùng là **end**

2. DAG with param

```
default_args = {
    "owner": "phamthixuanhien",
    "start_date": dt.datetime(2024, 4, 6, tzinfo=hcm_tz),
    "retries": 1,
    "retry_delay": dt.timedelta(minutes=5),
    "depend_on_past": True,
    "email": ["phamthixuanhien@iuh.edu.vn"]
}

def print_params(stt, email, sex):
    print(stt)
    print(email)
    print(sex)

with DAG(
    "dag_with_params",
    schedule=None,
    default_args=default_args,
    params={
        "stt": Param(1, type="integer", minimum=1),
        "email": 'phamthixuanhien@iuh.edu.vn',
        "sex": Param("Female", type="string", enum = ["Female", "Male"]),
    },
    catchup= False,
) as dag:

    start = EmptyOperator(task_id="start")

    print_string_task = PythonOperator(
        task_id="print_string",
        python_callable=print_params,
        op_kwargs={
            "stt": "{{ params.stt }}",
            "email": "{{ params.email }}",
            "sex": "{{ params.sex }}"
        },
    )

    end = EmptyOperator(task_id="end")

    # Flow
    start >> print_string_task >> end
```

params={...}: đây là một phần đặc biệt của đoạn mã. Đối tượng params được dùng để truyền các tham số tùy chỉnh vào DAG. 3 tham số params được định nghĩa sau:

- **"stt"**: tham số x với giá trị mặc định là 1, kiểu dữ liệu là số nguyên, và giá trị tối thiểu là 1
- **"email"**: giá trị mặc định là phamthixuanhien@iuh.edu.vn
- **"sex"**: tham số choice có giá trị mặc định là Female, kiểu dữ liệu là chuỗi, và chỉ chấp nhận các giá trị "Female" hoặc "Male"

Tạo một tác vụ sử dụng **PythonOperator**, loại tác vụ này cho phép thực thi một hàm Python trong quá trình chạy DAG

- **task_id**: đặt tên cho tác vụ
- **python_callable**=**print_params**: chỉ định hàm python (print_params) để chạy. Hàm này cần được định nghĩa ở nơi khác trong mã
- **op_kwargs**=**{...}**: truyền các tham số cho hàm **print_params**: Các giá trị lấy giá trị tương ứng từ đối tượng 'params' đã được định nghĩa trong DAG

3. DAG with Task_Group

```
with DAG(
    "dag_with_task_group",
    schedule_interval="0 0 * * *",
    default_args=default_args,
    catchup= False,
) as dag:
    start = EmptyOperator(task_id="start")

    @task_group()
    def task_group_1():
        EmptyOperator(task_id="do_something_1") >> [EmptyOperator(task_id="do_something_2"),
                                                       EmptyOperator(task_id="do_something_3")]

    end = EmptyOperator(task_id="end")

    # Flow
    start >> task_group_1() >> end
```

task_group_1(): được thực hiện bằng cách sử dụng **@task_group()**

Nhóm này bao gồm 3 emptyOperator với là công việc đầu tiên

task_id="do_something_1", sau đó 2 công việc song song là **task_id="do_something_2"** và **task_id="do_something_3"**

Giảng viên biên soạn: Phạm Thị Xuân Hiền

Trong đoạn mã này mô tả một luồng công việc đơn giản trong airflow với nhóm các công việc bao gồm công việc đơn lẻ và hai công việc song song. Điều này giúp tổ chức và quản lý các công việc phức tạp hơn một cách dễ dàng hơn.

4. DAG with PostgreSQL

Apache Airflow to run an SQL query on PostgreSQL use [PostgresOperator](#)

<pre>with DAG("dag_with_postgres", schedule_interval="0 0 * * *", default_args=default_args, catchup= False,) as dag: start = EmptyOperator(task_id="start") query = PostgresOperator(task_id = "postgres_query", postgres_conn_id="postgres_local", sql = 'select * from "QLSV"."SinhVien" ') end = EmptyOperator(task_id="end") # Flow start >> query >> end</pre>	<p><code>query = PostgresOperator(...)</code>: Tạo một task sử dụng <code>PostgresOperator</code>, dùng để thực hiện các truy vấn SQL trên cơ sở dữ liệu PostgreSQL. <code>task_id = "postgres_query"</code> đặt tên cho task, <code>postgres_conn_id="postgres_local"</code> định nghĩa kết nối đến cơ sở dữ liệu PostgreSQL (được cấu hình trước trong Airflow), và <code>sql = 'select * from "QLSV"."SinhVien" '</code> là câu lệnh SQL được thực hiện, trong trường hợp này là truy vấn tất cả dữ liệu từ bảng <code>SinhVien</code> trong schema <code>QLSV</code>.</p>
---	---

5. DAG with Bigquery

Apache Airflow to run an SQL query on Bigquery use [BigQueryExecuteQueryOperator](#)

<pre># Tạo DAG with DAG('bigquery_airflow_example', default_args=default_args, schedule_interval='@daily', # Chạy DAG hàng ngày catchup=False) as dag: start = EmptyOperator(task_id="start") # Tạo một task để thực hiện truy vấn trên BigQuery bq_query_task = BigQueryExecuteQueryOperator(task_id='bigquery_query', sql=''' SELECT * FROM `your-project.QLSV.SinhVien` WHERE MaBM = 'HTTT' LIMIT 100 ''', use_legacy_sql=False, # Sử dụng Standard SQL (không phải SQL Legacy của BigQuery) location='your-location', # Thay thế bằng vị trí của BQ dataset của bạn gcp_conn_id='your_gcp_connection', # Thay thế bằng ID của kết nối GCP trong Airflow) end = EmptyOperator(task_id="end", trigger_rule=TriggerRule.NONE_FAILED_MIN_ONE_SUCCESS) # Flow start >> bq_query_task >> end</pre>	
---	--

III. EXERCISE

Ex1: Làm lại 5 bài tập trên

lưu ý bài 4,5 cần cấu hình

Ex2:

Xây dựng một DAG trong Airflow để thực hiện các công việc sau:

- task 1: in ra dòng lệnh “Công việc 1”
- Task 2: viết một hàm xử lý dữ liệu `process_data()`, nội dung bên trong tùy ý.
- Task 3: viết một hàm lưu trữ dữ liệu `save_data()`, nội dung bên trong tùy ý

Yêu cầu:

- Đặt DAG ID là ‘MSSV_DAG2’ (MSSV của sinh viên)
- Định nghĩa các tham số trong DAG
 - Owner là họ tên của sinh viên
 - Không phụ thuộc vào các phiên bản trước đó (`depends_on_past: False`)
 - Thời điểm bắt đầu: ngày hiện tại theo múi giờ ASIA/HoChiMinh (giờ Việt Nam)
 - Số lần thử lại là 1 (`retries:1`)
 - Khoảng thời gian chờ giữa các lần thử lại là 2 phút
 - Thiết lập cho DAG chạy mỗi phút

Submission Regulations

- Thực hiện các bài tập sau, nhớ chụp màn hình kết quả vào word.
- Nộp file word và folder chứa code
- Quy tắc đặt tên: MSSV.zip