



CS 777: Big Data Analytics

Term Paper

APACHE AIRFLOW

Vaidehi Shah

Boston University, Metropolitan College
Spring 2023





INTRODUCTION

What is Apache Airflow?

Apache Airflow ^[1] is a powerful and versatile open-source platform. It is used for the efficient management and execution of data workflows, task scheduling, and orchestration of complex processes. At its core, Apache Airflow uses DAGs (Directed Acyclic Graphs).



Why Apache Airflow?

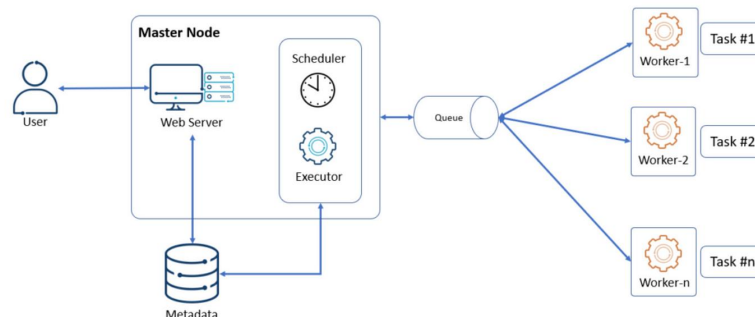
With the exponential growth of data processes, it addresses the need for efficient orchestration and automation of complex workflows and task scheduling. Airflow simplifies management by representing workflows as Directed Acyclic Graphs (DAGs), enabling automation and clear visualization of task dependencies. Some of its key features include:

- Flexibility : agnostic to the programming language or tools used in tasks
- Scalability : it has a distributed architecture
- Monitoring : provides a web-based UI that offers insights into the status of workflows, logs, and task histories; can also be integrated with external monitoring and alert systems



ARCHITECTURE

The architecture is based on a distributed and modular framework, which includes several core components:



1. **Metastore Database:** The Metastore Database, often powered by databases like PostgreSQL or MySQL, stores metadata related to tasks, task status, and workflow configurations. It keeps track of task history and maintains the state of the entire system.
2. **Scheduler:** The Scheduler is the brain of Apache Airflow. It constantly monitors all tasks and their dependencies and schedules when and how they should run. It ensures that tasks are executed in the correct order and takes into account task priorities, data dependencies, and resource availability.
3. **Executor:** Executors are responsible for executing tasks on worker nodes. Airflow supports various executors, such as the LocalExecutor, which runs tasks locally on the same machine as the Scheduler, and the CeleryExecutor, which uses a distributed task queue system like Apache Celery for task execution.
4. **Workers:** Workers are the compute resources that execute the actual tasks. They pull tasks from the message queue and execute them. Workers can be horizontally scaled to accommodate different workloads and ensure timely task execution.
5. **Web Interface:** The web-based UI provided by Airflow allows users to interact with the system, monitor task progress, and view historical data. It also provides a means to trigger, pause, and rerun tasks, making it a valuable tool for both developers and non-technical users.
6. **Message Queue:** Apache Airflow uses a message queue, such as Apache Kafka or RabbitMQ, to deliver tasks to workers. This decouples the scheduling of tasks from their execution, improving system reliability and enabling distributed task processing.
7. **DAGs (Directed Acyclic Graphs):** Workflows in Airflow are defined as DAGs, which are collections of tasks and their dependencies. These DAGs are written in Python code, providing a clear and programmable representation of workflow logic.



PROS & CONS OF AIRFLOW

Pros of Apache Airflow:

- **Flexibility and Extensibility:** Apache Airflow is highly flexible and extensible, allowing users to integrate with a wide range of technologies and customize workflows to their needs.
- **Dynamic Workflow Scheduling:** It optimizes task execution order based on dependencies and resource availability, reducing manual management and improving efficiency.
- **Monitoring and Alerting:** Provides a user-friendly interface and integrates with external tools for real-time monitoring, enhancing operational reliability.

Cons of Apache Airflow:

- **Learning Curve:** Airflow has a steep learning curve, particularly for newcomers to workflow orchestration.
- **Resource Consumption:** It can consume significant computational resources, necessitating careful resource management.
- **Complex Deployments:** Setting up high-availability or clustered configurations can be challenging and require expertise in distributed systems.

USE CASES

1. **ETL Processes:**

Apache Airflow is widely used for automating Extract, Transform, Load (ETL) processes, ensuring that data is efficiently extracted from source systems, transformed, and loaded into data warehouses or other destinations.

2. **Data Pipeline Automation:**

Organizations leverage Airflow to build and automate data pipelines that connect various data systems, enabling data ingestion, transformation, and delivery for analytics and reporting.

3. **Machine Learning Workflows:**

Airflow plays a crucial role in orchestrating machine learning workflows, streamlining data preparation, model training, hyperparameter tuning, and model deployment for data scientists and AI practitioners.



DEPLOYMENT ^[2]

Airflow installation using Docker (Docker provides isolated environment for running Airflow and its dependencies)

STEP 1

The first step is to install Docker, if you don't already have it on your system. You can go to the weblink ^[3] and install Docker for your specific system by following the instructions available there.

If you wish to check the successful installation of Docker and Docker-compose (will be done along with Docker) on your system, you can use the following commands:

```
>> docker --version
```

```
>> docker-compose --version
```

STEP 2

After this, it is recommended that you make a new folder in your system for airflow-docker, wherein you can download the docker-compose file describing all the services needed by airflow.

```
>> mkdir airflow-docker
```

```
>> cd airflow-docker
```

STEP 3

Here, you need to upload the docker-compose file, which has been made for us by the Airflow community. You can access the .yaml file ^[4]. To deploy Airflow on Docker Compose, you should fetch docker-compose.yaml:

```
>> curl -L -O 'https://airflow.apache.org/docs/apache-airflow/2.7.2/docker-compose.yaml'
```

STEP 4

Run the following command to start and run a set of Docker containers defined in a docker-compose.yml file in Docker Compose:

```
>> docker-compose up
```



STEP 5

You are now ready to access Airflow on one of your browsers: *Listening at: http://localhost:8080/login/*

The default username and password, both are set as "airflow".

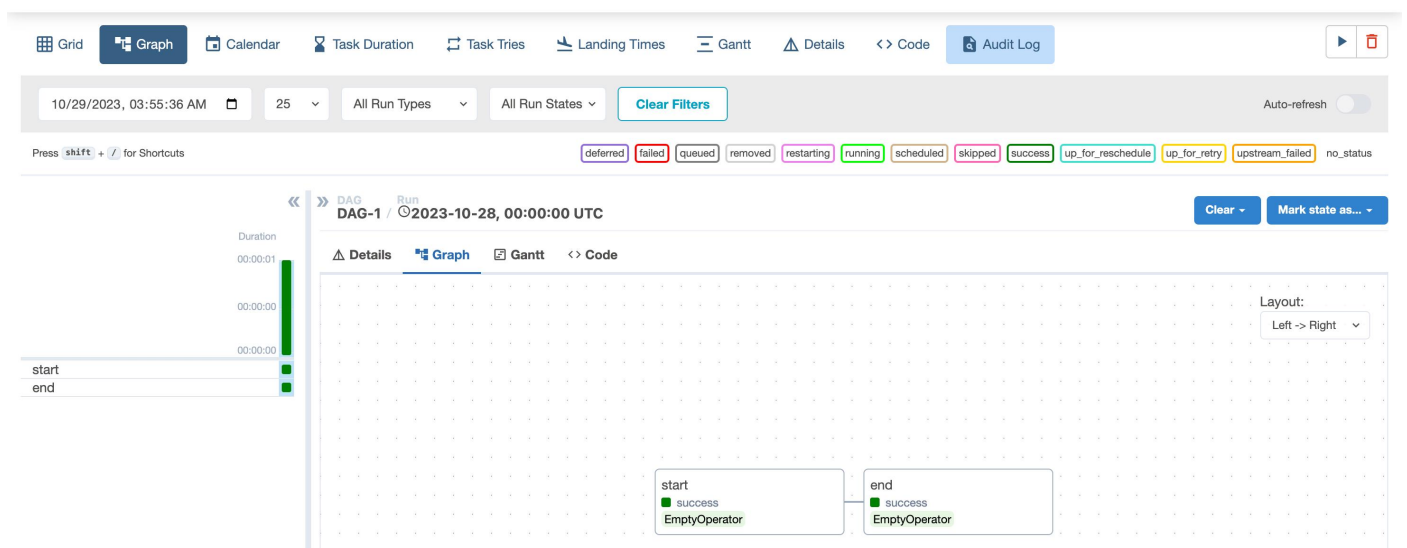
DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
DAG-1	airflow	1	@once	2023-10-28, 00:00:00		2
DAG-2	airflow	5	@once	2023-10-29, 00:13:41		5
DAG-3	airflow	2	None	2023-10-28, 19:33:41		3
dataset_consumes_1	airflow		Dataset		On s3://dag1/output_1.txt	
dataset_consumes_1_and_2	airflow		Dataset		0 of 2 datasets updated	
dataset_consumes_1_never_scheduled	airflow		Dataset		0 of 2 datasets updated	



WORKING WITH AIRFLOW

This is how I created my first DAG:

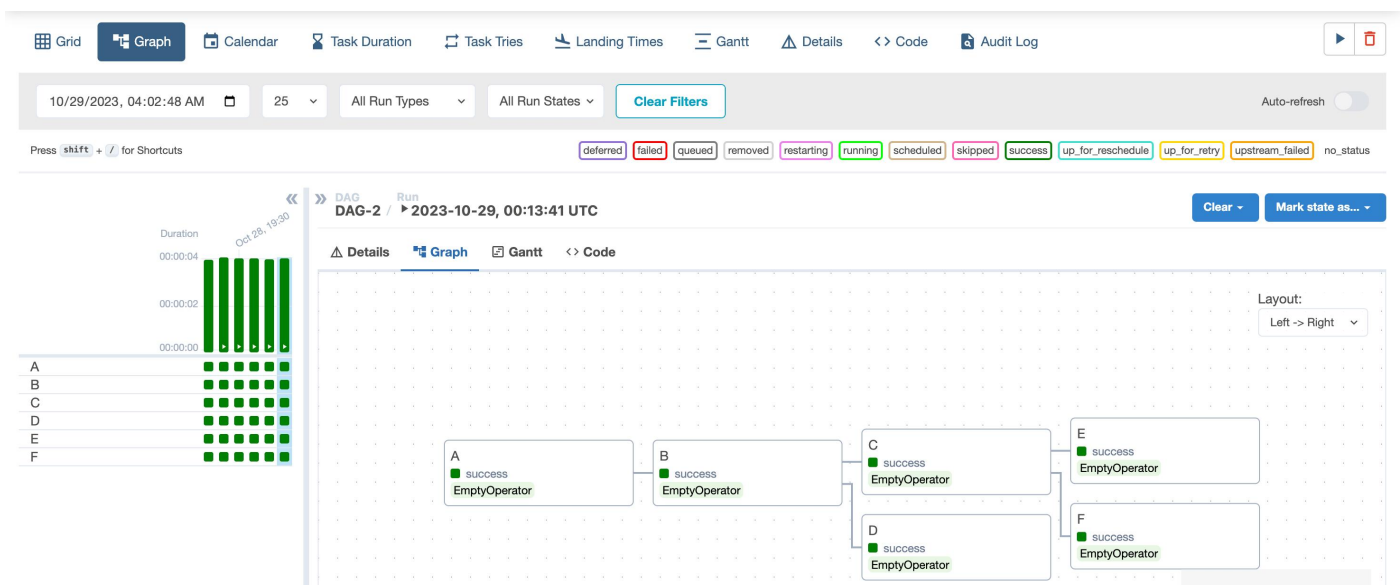
```
dag_1.py
1  # STEP: 1
2  ### importing modules
3  from airflow import DAG
4  from datetime import datetime, timedelta
5  from airflow.operators.dummy import DummyOperator
6
7  # STEP: 2
8  ### default args
9  default_args = {
10     'owner': 'airflow',
11     'depends_on_past': False,
12     'start_date': datetime(2023, 10, 28),
13     'retries': 0
14 }
15
16 # STEP: 3
17 ### instantiate DAG
18 dag = DAG(dag_id='DAG-1', default_args=default_args, catchup=False, schedule_interval='@once')
19
20 # STEP: 4
21 ### define tasks
22 start = DummyOperator(task_id='start', dag=dag)
23 end = DummyOperator(task_id='end', dag=dag)
24
25 # STEP: 5
26 ### define dependencies
27 start >> end
28
```





Demo code DAG-2:

```
dag_2.py
1 # STEP: 1
2 ### importing modules
3 from airflow import DAG
4 from datetime import datetime, timedelta
5 from airflow.operators.dummy import DummyOperator
6
7 # STEP: 2
8 ### default args
9 default_args = {
10     'owner': 'airflow',
11     'depends_on_past': False,
12     'start_date': datetime(2023, 10, 28),
13     'retries': 0
14 }
15
16 # STEP: 3
17 ### instantiate DAG
18 dag = DAG(dag_id='DAG-2', default_args=default_args, catchup=False, schedule_interval='@once')
19
20 # STEP: 4
21 ### define tasks
22 A = DummyOperator(task_id='A', dag=dag)
23 B = DummyOperator(task_id='B', dag=dag)
24 C = DummyOperator(task_id='C', dag=dag)
25 D = DummyOperator(task_id='D', dag=dag)
26 E = DummyOperator(task_id='E', dag=dag)
27 F = DummyOperator(task_id='F', dag=dag)
28
29 # STEP: 5
30 ### define dependencies
31 A >> B
32 B >> C
33 B >> D
34 C >> [E, F]
```





REFERENCES

- [1] <https://airflow.apache.org/docs/apache-airflow/stable/index.html>
- [2] https://github.com/vaidehi-nshah/airflow_demo
- [3] <https://docs.docker.com/desktop/>
- [4] <https://airflow.apache.org/docs/apache-airflow/stable/howto/docker-compose/index.html>
- [5] https://youtu.be/K9AnJ9_ZAXE?feature=shared
- [6] <https://youtu.be/aTaytcxy2Ck?feature=shared>