Airflow

# What is Apache Airflow

**Apache Airflow** is an open-source platform to Author, Schedule and Monitor workflows.

It was created at Airbnb and currently is a part of **Apache** Software Foundation.

**Airflow** helps you to create workflows using Python programming language and these workflows can be scheduled and monitored easily with it.

# Some use cases

Automate training, testing and deploying a machine learning model

Ingesting data from multiple REST APIs

Coordinating Extraction, Transformation and Loading (ETL) or Extraction, Loading and Transformation (ELT) Processes Across an Enterprise Data Lake

As we can see, one of the main features of **Airflow** is its flexibility: it can be used for many different data workflow scenarios.

Due to this aspect and its rich feature set, it has gained significant traction over the years.
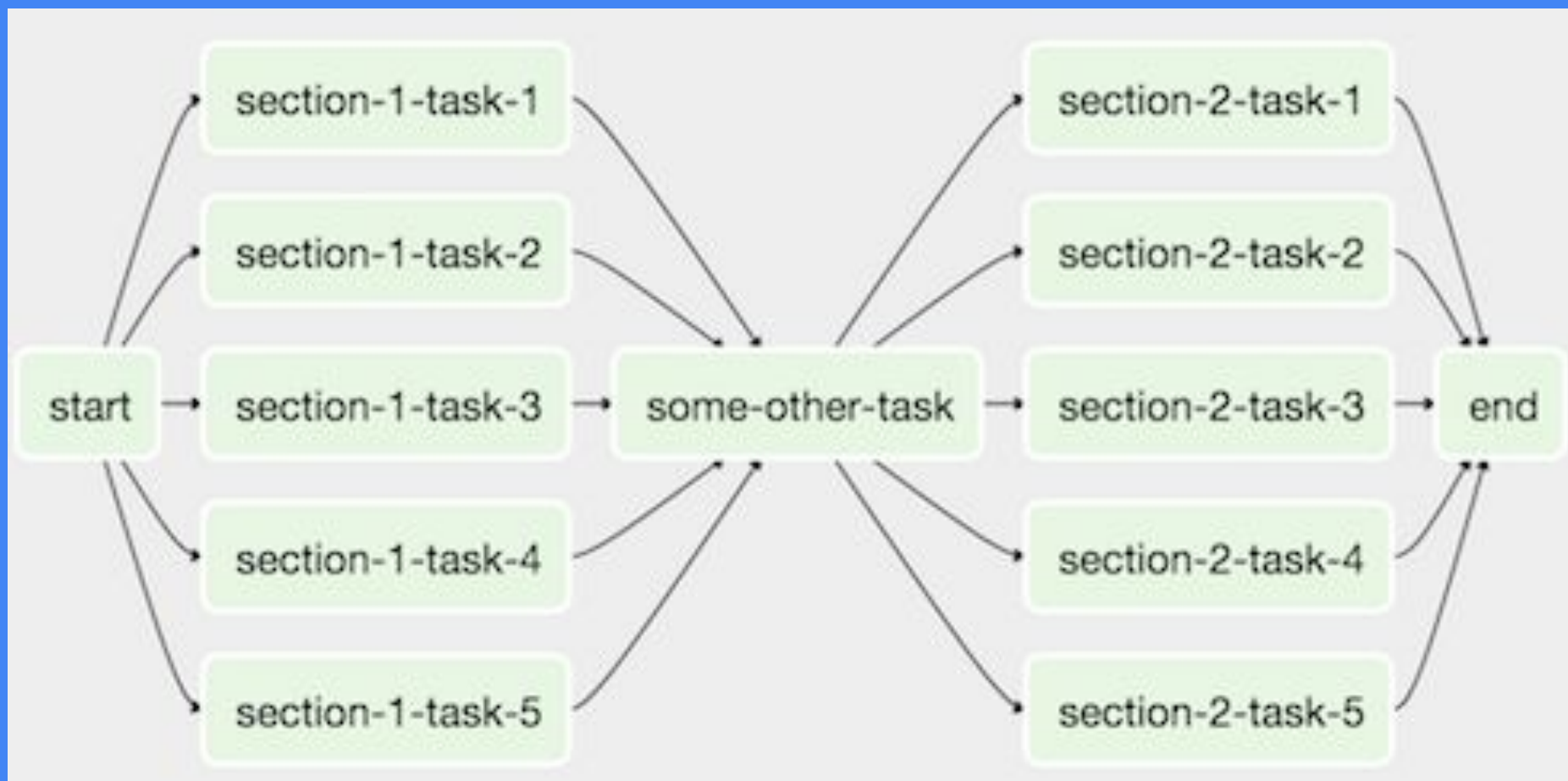
It has been battle-tested by many companies, from startups to Fortune 500 enterprises.

Some examples include Spotify, Twitter, Walmart, Slack, Robinhood, Reddit, PayPal, Lyft, and of course, Airbnb.

# Apache DAGs

**Apache Airflow** works with the concept of Directed Acyclic Graphs (DAGs), which are a powerful way of defining dependencies across different types of tasks.

In **Apache Airflow**, DAGs are developed in **Python**, which unlocks many interesting features from software engineering: modularity, reusability, readability, among others.

```python
26
27  t_download = BashOperator(
28      task_id='download_today',
29      bash_command='python3 ' + scriptpath + 'downloadtoday.py "{{ execution_date }}"',
30      retries=3,
31      dag=dag)
32
33  t_extract = BashOperator(
34      task_id='extract_downloads',
35      bash_command='python3 ' + scriptpath + 'extractdownloads.py "{{ execution_date }}"'
        ,
36      retries=3,
37      dag=dag)
38
39  t_singerconfig = BashOperator(
40      task_id='singer_config',
41      bash_command='python3 ' + scriptpath + 'generatesingerconfig.py "{{ execution_date
          }}"',
42      retries=3,
43      dag=dag)
44
45  t_singerpush = BashOperator(
46      task_id='singer_push',
47      bash_command='tap-csv -c ' + base_path + "{{ execution_date.strftime('%Y-%m-%d') }}
          " + '/' + 'csv-config.json | target-stitch -c ~/Stitch/configs/
          stitch_mbi_config.json',
48      retries=3,
49      dag=dag)
50
51  t_cleanup = BashOperator(
52      task_id='cleanup',
53      bash_command='python3 ' + scriptpath + 'cleanup.py "{{ execution_date }}"',
54      retries=3,
55      dag=dag)
56
57  t_cleanup.set_upstream(t_singerpush)
58  t_singerpush.set_upstream(t_singerconfig)
59  t_singerconfig.set_upstream(t_extract)
60  t_extract.set_upstream(t_download)
```

# Apache Sensors, Hooks, Operator

**Sensors, Hooks, and Operators** are the main building blocks of Apache Airflow. They provide
an easy and flexible way of connecting to multiple external resources. Want to integrate with Amazon S3? Maybe you also want to add Azure Container Instances to the picture in order to run some short-lived Docker containers? Perhaps running a batch workload in an Apache Spark or Databricks cluster? Or maybe just executing some basic Python code to connect with REST APIs? Airflow ships with multiple operators, hooks, and sensors out of the box, which allow for easy integration with these resources, and many more, such as DockerOperator, BashOperator, HiveOperator, JDBCOperator — the list goes on. You can also build upon one of the standard operators and create your own. Or you can simply write your own operators, hooks, and sensors from scratch.

Airflow UI

The UI allows for quick and easy monitoring and management of your Airflow instance.

Detailed logs also make it easier to debug your DAGs.

**On** DAG: example_dag

✳ Graph View    ♦ Tree View    �ⅈⅈ Task Duration    📑 Task Tries    ✈ Landing Times    ☰ Gantt    ☰ Details    ⚡ Code    ⟳ Refresh    ⊗ Delete

**running** Base date: `2019-10-22 12:06:01`    Number of runs: `5 ▾`    Run: `scheduled__2019-10-22T12:06:00+00:00 ▾`    Layout: `Left->Right ▾`    Go

PythonProgramOperator

# DAGs

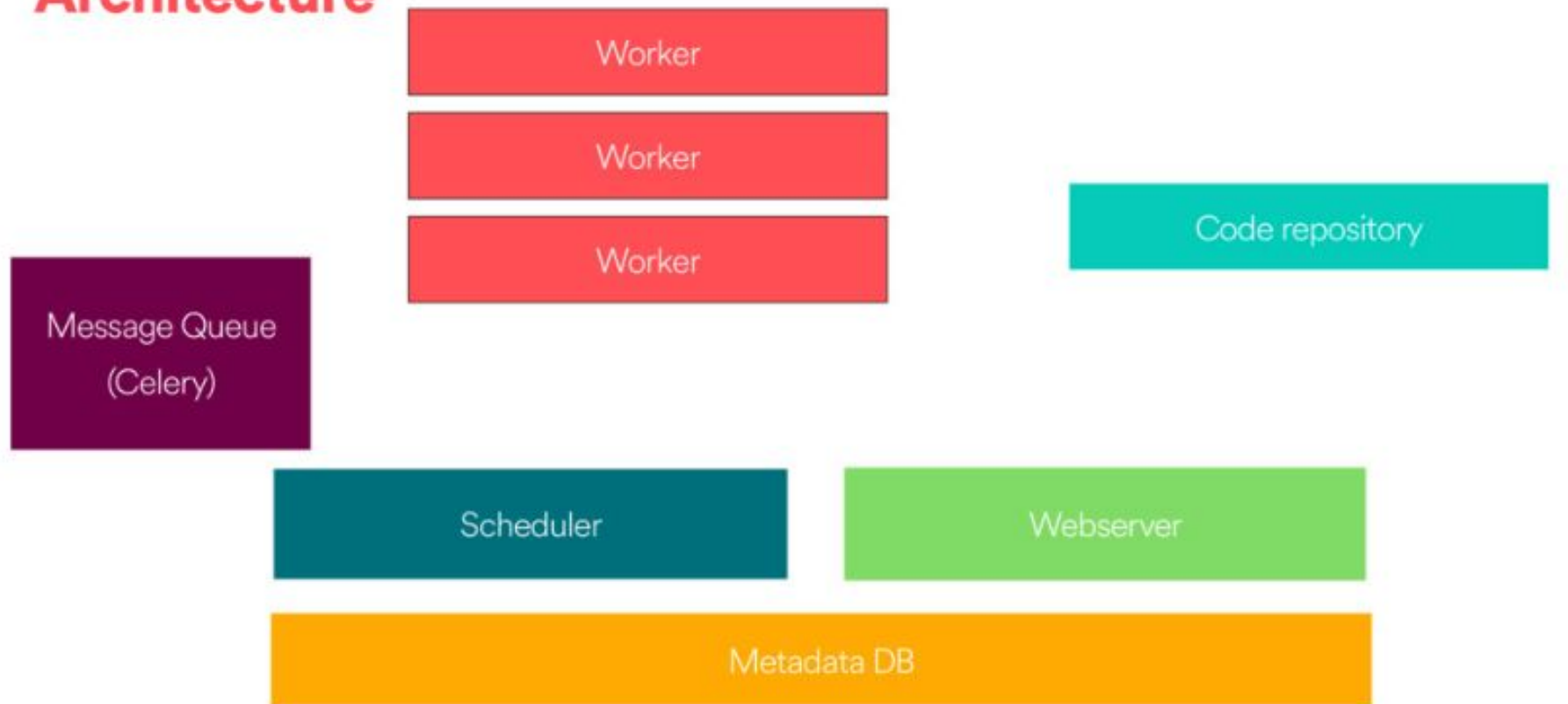| | | DAG | Schedule | Owner | Recent Tasks ⓘ | Last Run ⓘ | DAG Runs ⓘ | Links |
|---|---|---|---|---|---|---|---|---|
| | On | example_bash_operator | 0 0 * * * | Airflow | | | | |
| | On | example_gcs_to_sftp | None | airflow | | | | |
| | On | example_trigger_controller_dag | @once | airflow | | | | |
| | On | example_trigger_target_dag | None | Airflow | | | | |
| | On | test_backfill_pooled_task_dag | 1 day, 0:00:00 | airflow | | | | |
| | On | test_default_impersonation | 1 day, 0:00:00 | airflow | | | | |
| | On | test_task_view_type_check | 1 day, 0:00:00 | airflow | | | | |

Showing 1 to 7 of 7 entries

« ‹ 1 › »

Hide Paused DAGs

# Airflow Components

# Architecture

Worker

Worker

Worker

Code repository

Message Queue
(Celery)

Scheduler

Webserver

Metadata DB

# Web Server

The GUI. This is under the hood a Flask app where you can track the status of your jobs and read logs from a remote file store

# Scheduler

This component is responsible for scheduling jobs. This is a multithreaded Python process that uses the DAGb object to decide what tasks need to be run, when and where.

The task state is retrieved and updated from the database accordingly. The web server then uses these saved states to display job information.

# Executor

The mechanism that gets the tasks done.

# Metadata Database

- Powers how the other components interact

- Stores the Airflow states

- All processes read and write from here

# Questions?