

Naya College

Apache Airflow



Apache Airflow



V 1.0.0

Agenda

Part 1

- History
- Introduction to Airflow
- Why Airflow?
- Airflow Architecture
- Airflow Components
- User Interface

Part 2

- Airflow Workflow in Details
- Creating an Airflow Workflow Methods

Part 3

- Running a Job Using Airflow CLI
- Creating & Running DAGs in Airflow Workflow
- Viewing & Troubleshooting DAGs Airflow Workflows



History

Airbnb is a fast growing, data informed company.

The data teams and data volume are growing quickly, and accordingly, so does the complexity of the challenges we take on.

Airbnb growing workforce of data engineers, data scientists and analysts are using **Airflow**, a platform we built to allow us to move fast, keep our momentum as we author, monitor and retrofit **data pipelines**.

Today, **Airflow** is **open sourcing** and **sharing** Airflow.



Introduction to Airflow

Airflow is a platform to programmatically author, Schedule and monitor workflows.

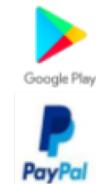
Use Airflow to author workflows as Directed Acyclic Graphs (DAGs) of tasks. The Airflow scheduler executes your tasks on an array of workers while following the specified dependencies.

Rich command line utilities make performing complex surgeries on DAGs a snap.

The rich user interface makes it easy to visualize pipelines running in production, monitor progress, and troubleshoot issues when needed. When workflows are defined as code, they become more maintainable, versionable, testable, and collaborative.

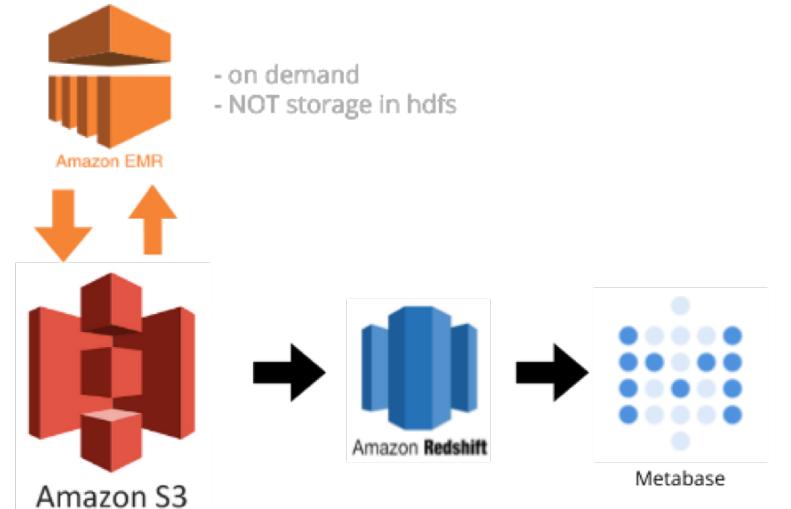
Introduction to Airflow

Data Sources



mongoDB

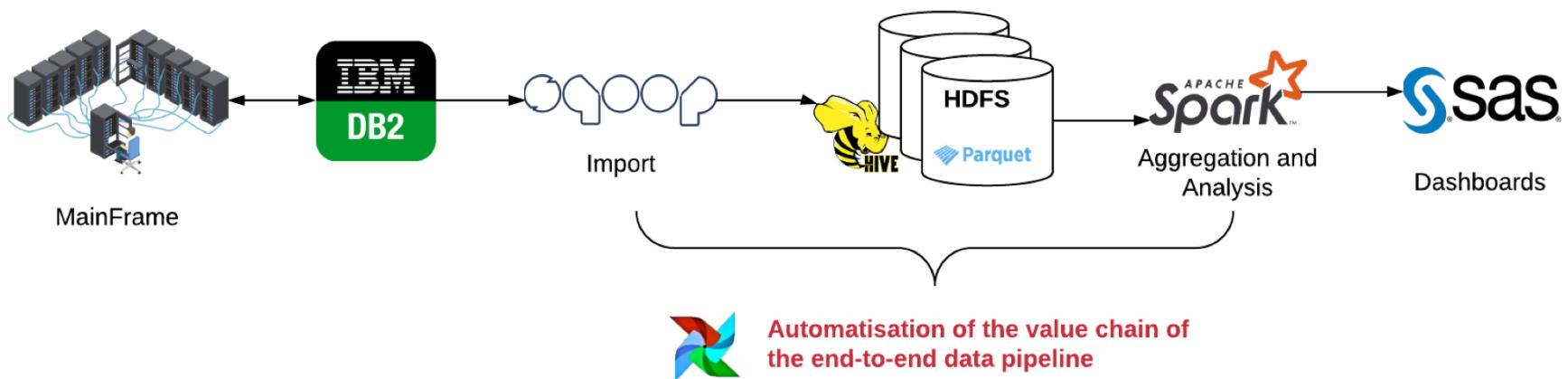
stripe



Introduction to Airflow

Airflow is a platform to programmatically author,
Schedule and monitor workflows.

Big Data Architecture



Why Airflow?

- One Place for schedule and data pipeline
- Python code base
- Callable events
- Xcoms
- Trigger rules
- Simply UI
- Powerful CLI
- Queues & Pools
- Growing community

The screenshot shows the Apache Airflow web interface with a dark blue header bar. The header includes the Airflow logo, navigation links for DAGs, Data Profiling, Browse, Admin, Docs, and About, and a timestamp: 2018-09-07 22:14:10 UTC.

The main content area is titled "DAGs". It features a table with the following columns: DAG, Schedule, Owner, Recent Tasks, Last Run, DAG Runs, and Links. The table lists five DAG entries:

DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
example_bash_operator	@0****	airflow	1	2018-09-06 00:00	1	0
example_branch_dop_operator_v3	@1*****	airflow	1, 1	2018-09-05 00:56	1	0
example_branch_operator	@daily	airflow	1	2018-09-06 00:00	1	0
example_xcom	@once	airflow	1	2018-09-05 00:00	1	0
latest_only	@00:00	Airflow	1	2018-09-07 16:00	1	0

Below the table, there is a search bar labeled "Search:" and a pagination control with the number "1" highlighted. A small note at the bottom says "Showing 1 to 5 of 5 entries".

Components

- **Task** –a defined unit of work (these are called operators in Airflow)
- **DAG** –Directed acyclic graph, a set of tasks with explicit execution order, beginning, and end

Off DAG: imap_example_dag

Graph View Tree View Task Duration Task Tries Landing Times Gantt Details Code Refresh Delete

running Base date: 2019-08-04 00:00:01 Number of runs: 25 Run: scheduled__2019-08-04T00:00:00+00:00 Layout: Left->Right Go

DummyOperator ExtendedPythonOperator IMAPAttachmentOperator

```
graph LR; start((start)) --> get_email_attachment[get_email_attachment]; get_email_attachment --> upload_attachment_to_s3[upload_attachment_to_s3]; upload_attachment_to_s3 --> remove_local_attachment[remove_local_attachment]; remove_local_attachment --> end((end))
```

Components

- **Hooks** – Interfaces to external platforms and databases like Hive, S3, MySQL, Postgres, HDFS, and Pig
- **Airflow pools** – can be used to limit the execution parallelism on arbitrary sets of tasks.
- **Variables** – Generic way to store and retrieve arbitrary content or settings as a simple key value store within Airflow

Variables

Choose file No file chosen Import Variables

List (2) Create Add Filter With selected Search: key, val, is_encrypt

	Key	Val
<input type="checkbox"/>	v_name	#####THIS IS A SOME VALUE#####
<input type="checkbox"/>	v_st_dt	2019-01-02

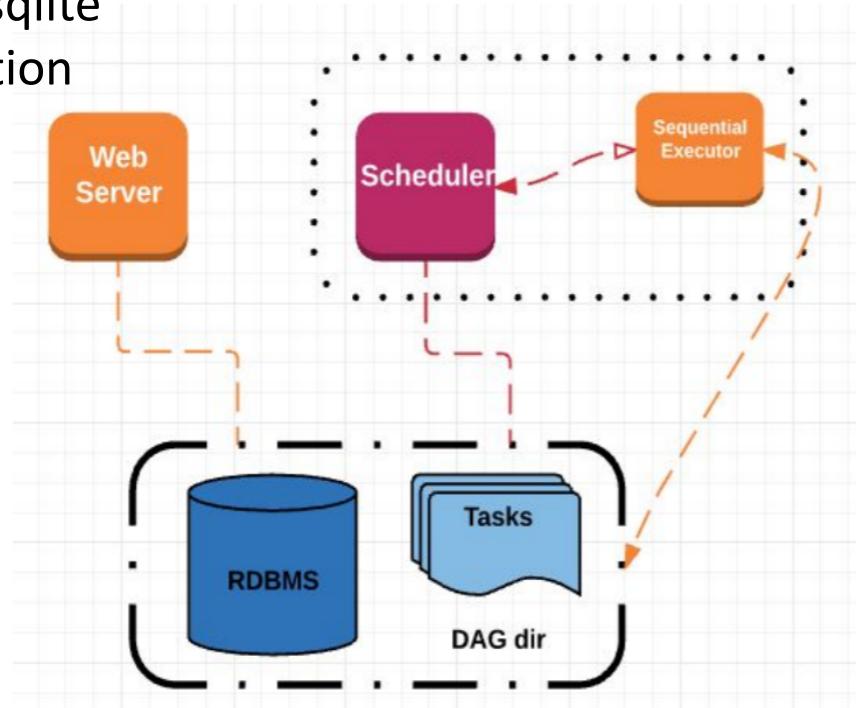
Components

- **Queues** – When using the CeleryExecutor, the celery queues that tasks are sent to can be specified. queue is an attribute of BaseOperator.
- **SubDAGs** –Defining a function that returns a DAG object is a nice design pattern when using Airflow.
- **SLAs** –Service Level Agreements, or time by which a task or DAG should have succeeded, can be set at a task level as a timedelta
- **Trigger Rules** –Though the normal workflow behavior is to trigger tasks when all their directly upstream tasks have succeeded, Airflow allows for more complex dependency settings.

Architecture

Sequential

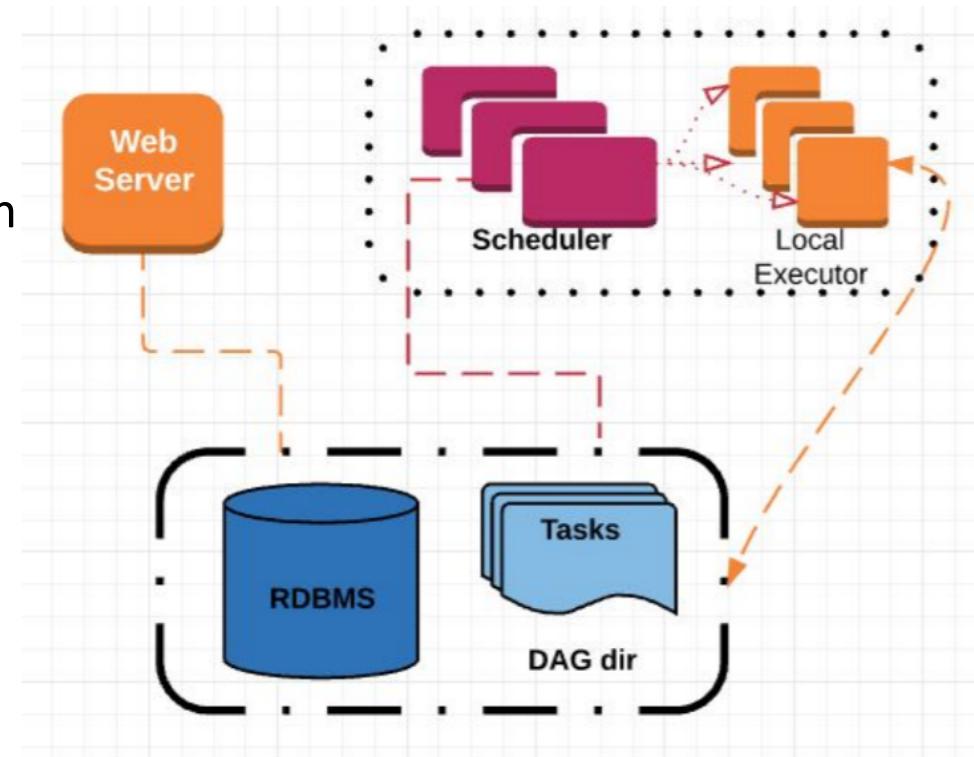
- Process one task at time (not parallelism)
- Default mode
- Minimum setup – works with sqlite
- Not recommended for production



Architecture

Local Executor

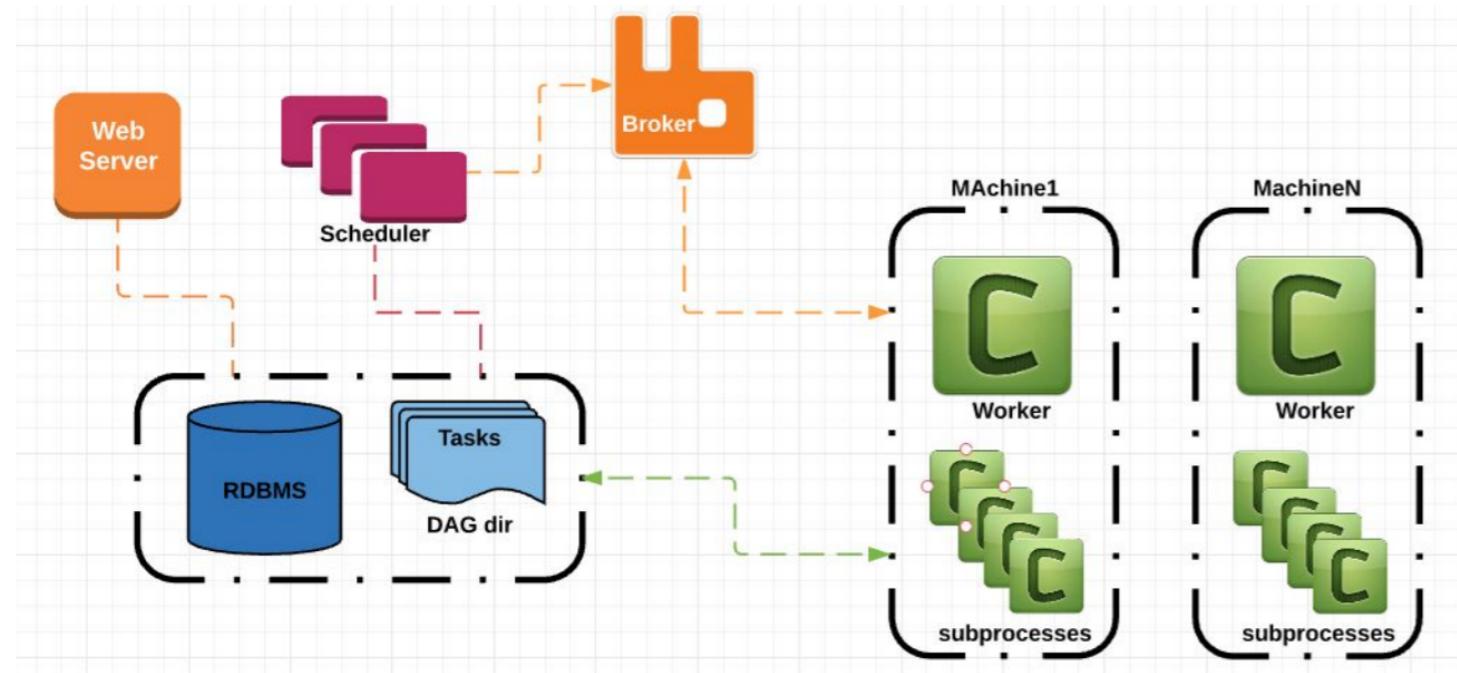
- Spawned by scheduler processes
- Scales vertical
- Production grade
- Doesn't need broker
- Suitable for production when there's not so many DAGs



Architecture

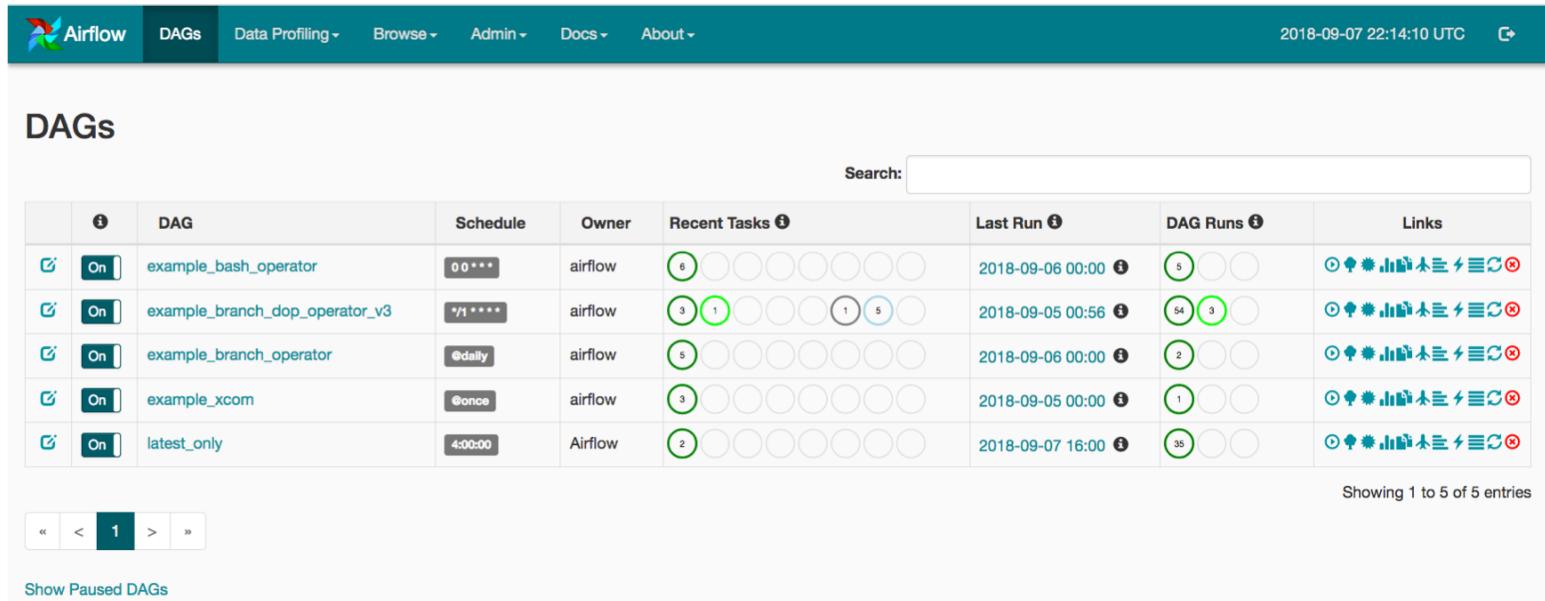
Celery Execute

- Vertical and Horizontal scalable
- Can be monitored (via Flower) Production grade
- Support Pools and Queues



User Interface

List of DAGs



The screenshot shows the Apache Airflow user interface for managing Data Automation Groups (DAGs). The top navigation bar includes links for Airflow, DAGs, Data Profiling, Browse, Admin, Docs, and About, along with the current timestamp (2018-09-07 22:14:10 UTC) and a refresh icon.

The main content area is titled "DAGs" and features a search bar labeled "Search:". Below the search bar is a table listing five DAG entries:

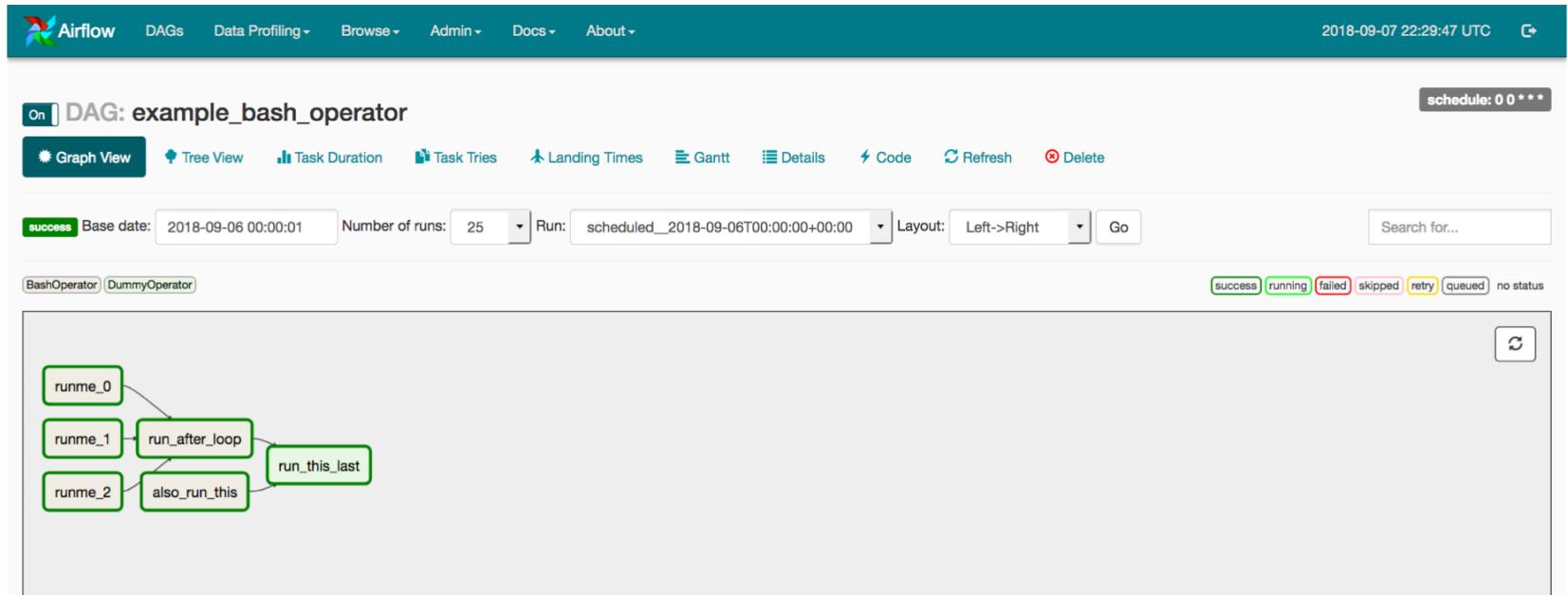
	DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
<input checked="" type="checkbox"/>	On example_bash_operator	<code>0 * * * *</code>	airflow	6 ○ ○ ○ ○ ○ ○	2018-09-06 00:00 i	5 ○ ○	① ⌚ ✖ 🕒 🕒 🕒 🕒
<input checked="" type="checkbox"/>	On example_branch_dop_operator_v3	<code>*1 * * * *</code>	airflow	3 1 ○ ○ ○ ○ ○	2018-09-05 00:56 i	54 3 ○	① ⌚ ✖ 🕒 🕒 🕒 🕒
<input checked="" type="checkbox"/>	On example_branch_operator	<code>@daily</code>	airflow	5 ○ ○ ○ ○ ○	2018-09-06 00:00 i	2 ○ ○	① ⌚ ✖ 🕒 🕒 🕒 🕒
<input checked="" type="checkbox"/>	On example_xcom	<code>@once</code>	airflow	3 ○ ○ ○ ○ ○	2018-09-05 00:00 i	1 ○ ○	① ⌚ ✖ 🕒 🕒 🕒 🕒
<input checked="" type="checkbox"/>	On latest_only	<code>4:00:00</code>	Airflow	2 ○ ○ ○ ○ ○	2018-09-07 16:00 i	35 ○ ○	① ⌚ ✖ 🕒 🕒 🕒 🕒

At the bottom left, there are navigation buttons for page 1 of 1. At the bottom right, it says "Showing 1 to 5 of 5 entries".

[Show Paused DAGs](#)

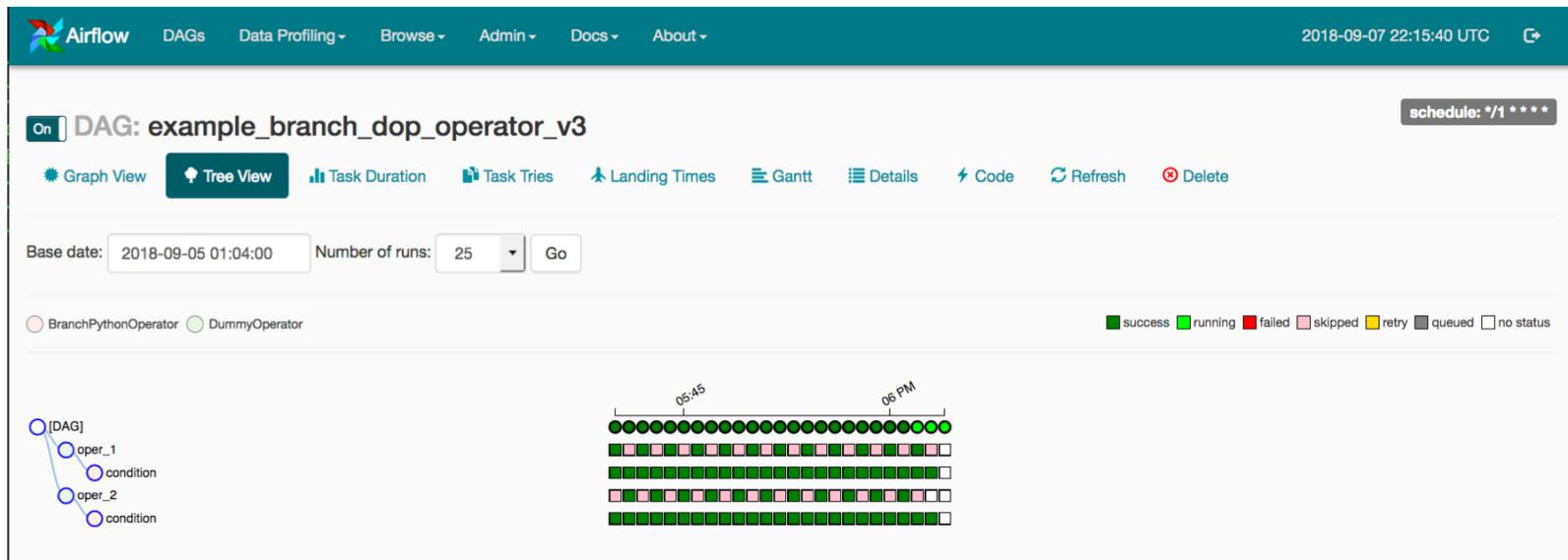
User Interface

Graph View



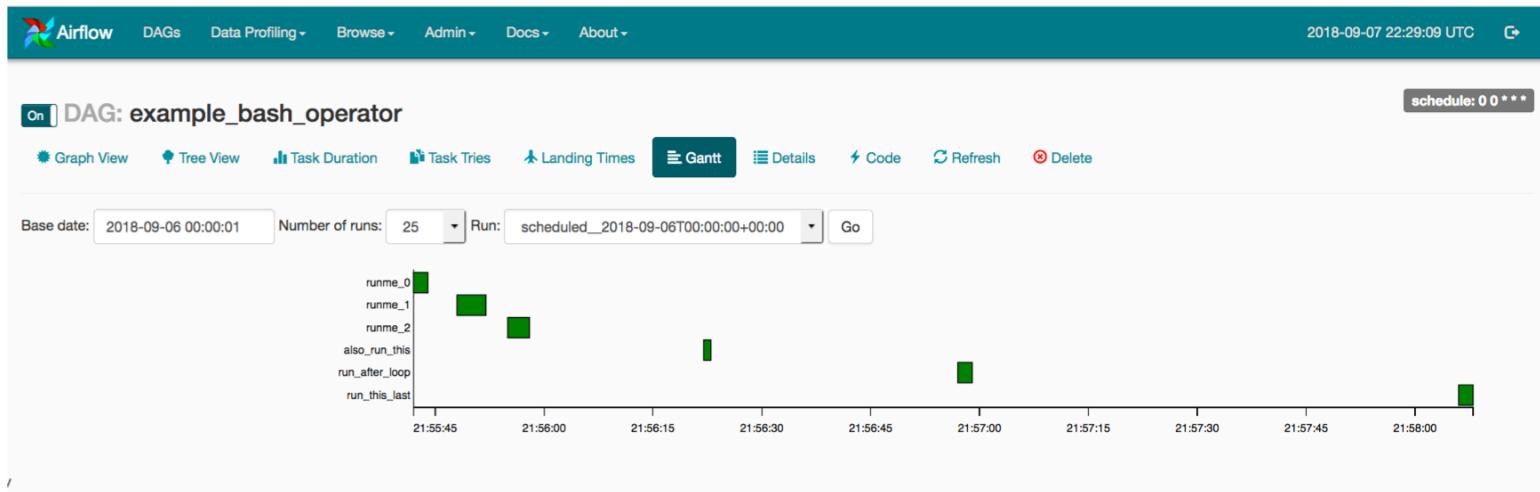
User Interface

Tree View



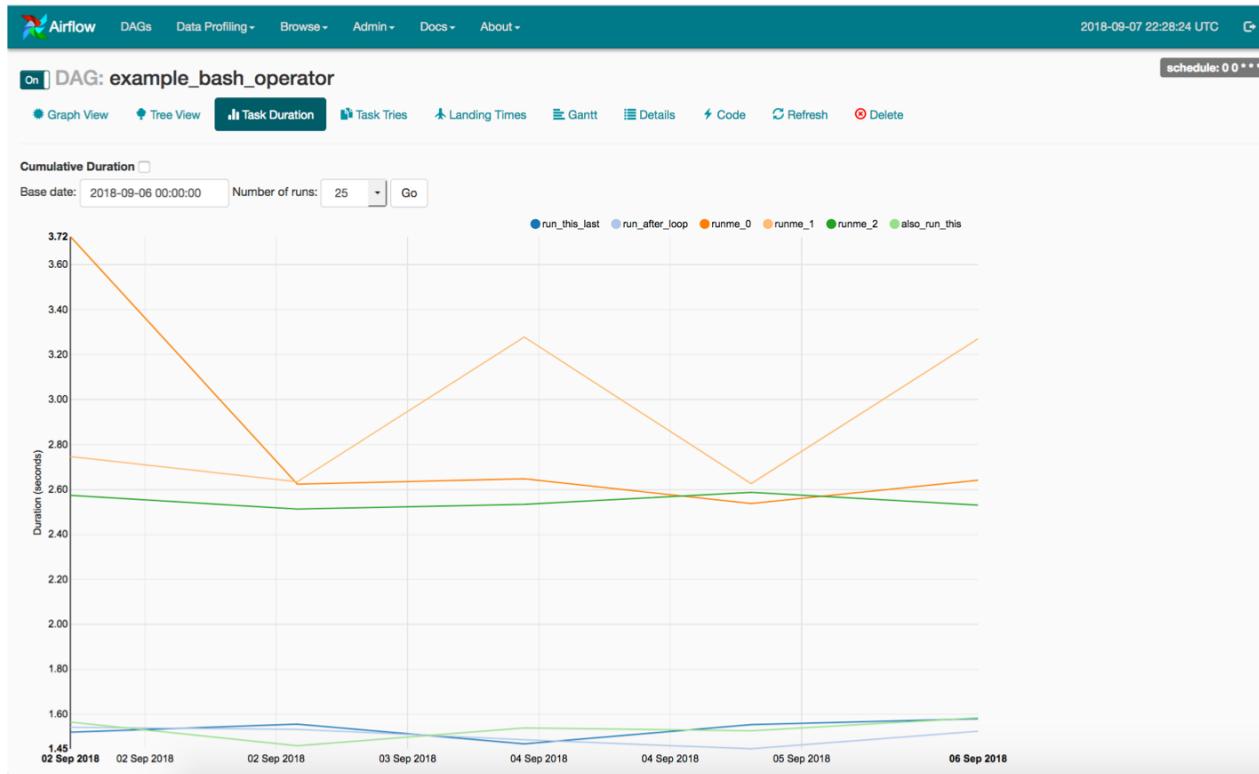
User Interface

Gantt Chart



User Interface

Task Duration



Airflow

Tutorial