

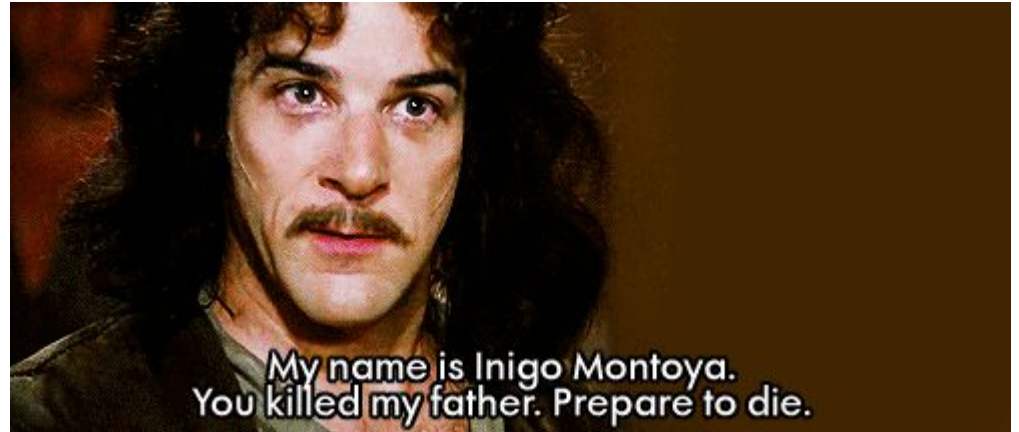
Airflow and DBT

Day one: Airflow



Introductions

- Polite Greeting
- Your Name
- Relevant Personal Link
- Clear Expectations



What is Airflow

Apache Airflow is an open-source workflow management platform for creation, scheduling, & management of data-engineering pipelines

- Implement programs from any language
- Workflows are written in Python
- Implements Workflows as Directed Acyclic Graphs (DAGs)
- Workflows are accessed via code, command line, or web-interface/REST API

What is dbt

Data Build Tool (dbt) is a development framework combining modular SQL with software engineering best-practices to enable reliable and fast data transformation

- Write custom business logic using SQL
- Build data pipelines
- Automate data quality testing
- Deliver data with documentation side-by-side with code

Day one: The plan

- 10.00 Airflow Basics:
 - Use case
 - Airflow in a nutshell
- 10.45 Short Break
- 11.00 Hands-on: Building a basic DAG
- 12.00 Lunch I'm sure
- 13.00 Airflow: Let's dive deeper
 - Internal architecture
 - XCOM
 - DAGs
 - Sensors
- 14.00 Short Break
- 14.15 Custom Operator: Python Code
- 16.00 Airflow Use Case: Centralized Pipelines
- 16.30 End

Use-case: Transparent pipelines with Airflow

We had a large, unorganized dataset of questionable quality. We were asked to fix the dataset to improve quality and enable the extraction of credible market insights for the company. We identified goals including:

- Extracting data from the central data team (AWS Redshift)
- Deduplication
- Cleaning
- Filling in missing data using business logic
- Filling in missing data or enhancing data using Machine Learning models
- Loading the final enhanced dataset into a PostgreSQL database for the webapp

How could we create a clear, transparent, and manageable pipeline? We chose Airflow

Use-case: Transparent pipelines with Airflow (2)

Airflow allowed us to:

- Build transparency into the workflow
- Optimize the development workflow
- Easily re-run failed Tasks
- Deliver insight into when Tasks ran, for how long, etc
- Integrate services as needed (Python, SQL, Redshift, Kubernetes, Lambdas)

And because Airflow is Python, it was fairly easy to get started on migrating our pipeline to a clean setup with CI/CD, including unit tests.

Airflow in a nutshell

Apache Airflow is an open-source workflow management platform for creation, scheduling, & management of data-engineering pipelines

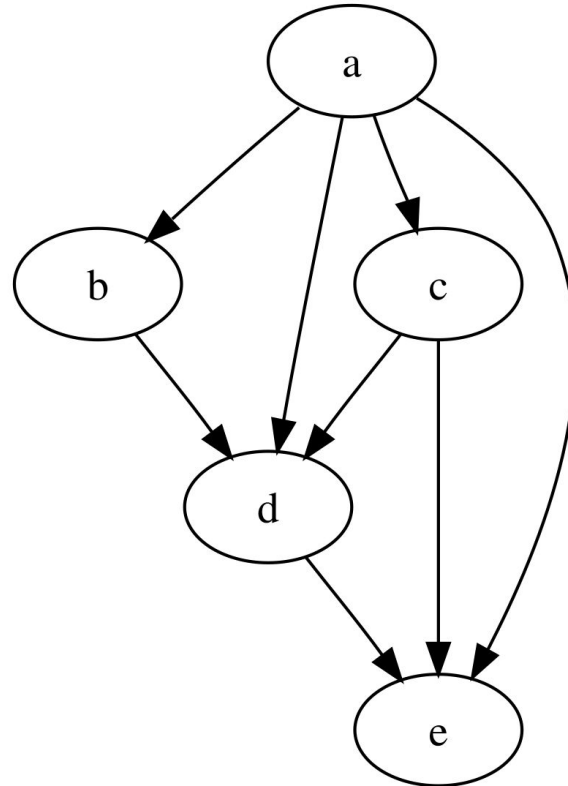
- DAGs
- Tasks
- UI features
- Jinja templating

What are DAGs

Directed Acyclic Graph (DAG)

A one-directional flow-chart of tasks

- Tasks
- Task dependencies
- Task details (Name, start date, owner, etc.)



What are Tasks?

Tasks are the basic executable unit in Airflow.

Tasks are organized into **DAGs**, with upstream and downstream dependencies (the flowchart / **DAG**) and are executable.

There are three types:

- **Operators**: pre-defined, reusable Tasks
- **Sensors**: a Task type that waits until a condition is met before executing
- **TaskFlow decorators**: custom python code (@task)

Operators: Pre-defined, reusable Tasks

- BashOperator: Executes a Bash script
- PythonOperator: Executes python code
- EmailOperator: Sends an email
- KubernetesPodOperator: Interacts with Kubernetes Pods

And many, many more. Some are also ecosystem dependent.

```
1 ▼ kubernetes_min_pod = GKStartPodOperator(  
2     # The ID specified for the task.  
3     task_id="pod-ex-minimum",  
4     # Name of task you want to run, used to generate Pod ID.  
5     name="pod-ex-minimum",  
6     project_id=PROJECT_ID,  
7     location=CLUSTER_REGION,  
8     cluster_name=CLUSTER_NAME,  
9     # Entrypoint of the container, if not specified the Docker container's  
10    # entrypoint is used. The cmds parameter is templated.  
11    cmds=["echo"],  
12    # The namespace to run within Kubernetes, default namespace is  
13    # `default`.  
14    namespace="default",  
15    # Docker image specified. Defaults to hub.docker.com, but any fully  
16    # qualified URLs will point to a custom repository. Supports private  
17    # gcr.io images if the Composer Environment is under the same  
18    # project-id as the gcr.io images and the service account that Composer  
19    # uses has permission to access the Google Container Registry
```

```
1 ▼ email_notification = EmailOperator(  
2     task_id="email_notification",  
3     to="noreply@astronomer.io",  
4     subject="dag completed",  
5     html_content="the dag has finished",  
6 )  
7
```

Sensors: They sense things.

Sensors are Operators designed to execute when triggered by an event. Sensors can be time-based, wait for a file, or an external event, but they wait for a condition to be met and then execute so downstream Tasks can run.

Examples:

- **DayTimeSensor**: Waits for specific day and time
- **ExternalTaskSensor**: Triggers when another Task finishes
- **SQLSensor**: Waits for data to arrive

Sensors have one flaw, and we will discuss that this afternoon.

Taskflow decorators - Creating Tasks with API

TaskFlow decorators such as `@task` are programming patterns that add functionality to an object. In this case, they turn functions into tasks with an ID, facilitate passing data between Tasks, and define dependencies. Decorators create similar functionality to and are interoperable with classical Operators, including `PythonOperator`.

TaskFlow is newer and therefore preferable over building the same functionality in `PythonOperator`.

UI features

The Airflow UI offers a feature-rich interface for:

- Monitoring Workflows and Tasks
- Visualization for Workflows
- Visualizing connections
 - These are also a one-stop for auth information for DAGs and Tasks
- Inspect DAGs and code from the UI (easier debugging)
- One-stop for Gantt charts, logs, timing, re-running, etc.

Jinja Templating

Jinja is a lightweight templating engine for Python that embeds variables, logic, and control structures to dynamically generate text. It enables flexible configuration and parameterization, especially useful for defining dynamic content like SQL queries or configuration files.

Jinja improves DAGs by dynamically injecting task parameters, SQL code, and runtime values (like execution date) at runtime. This improves reusability, reduces hardcoding, and supports complex automation logic across time-based data pipelines.

Airflow use-cases: Why?

- **Schedule workflows**: Automates recurring tasks like daily data jobs or reports.
- **Visual DAGs**: Clear, graphical view of task dependencies and execution status.
- **Modular Python code**: Workflows are written in Python, making them readable and reusable. Also easier to use than Infrastructure as Code variants like AWS step functions.
- **Agnostic**: Airflow is an Apache tool and can be used locally or with any cloud platform.
- **Handles failures**: Built-in retry, alerting, and logging for robust pipelines.
- **Track and debug**: Easy to inspect past runs, logs, and task outputs.
- **Extensible**: Supports plugins, custom operators, and integrates with many tools.
- **Scalable**: Run small local jobs or scale out with distributed executors (Celery, Kubernetes).
- **Templating & parameters**: Dynamically generate SQL, API calls, or filenames with runtime context.

Short Break

Morning exercise

Let's get started with Airflow! On the laptops, you'll find that Docker, WSL and Airflow are installed. Our first order of business:

- Open Docker and run the container for Airflow.

This is the equivalent of a fresh install — the UI is there, but we haven't created any DAGs yet.

Let's create a DAG

```
main.md x from airflow import DAG
1 from airflow import DAG
2 from airflow.operators.bash import BashOperator
3 from airflow.utils.dates import days_ago
4 from datetime import timedelta
5
6 with DAG(
7     dag_id="simple_bash_graph",
8     description="Demo DAG showing task dependencies with BashO
9     schedule_interval=None, # Manual trigger only
0     start_date=days_ago(1),
1     catchup=False,
2     default_args={
3         "retries": 0,
4         "retry_delay": timedelta(minutes=1),
5     },
6     tags=["example"],
7 ) as dag:
8
9     bash1 = BashOperator(
0         task_id="bash1",
```

```
bash1 = BashOperator(
    task_id="bash1",
    bash_command='echo "Running bash1 - {{ run_id }}"',
)

bash2 = BashOperator(
    task_id="bash2",
    bash_command='echo "Running bash2 - {{ run_id }}"',
)

bash3 = BashOperator(
    task_id="bash3",
    bash_command='echo "Running bash3 - {{ run_id }}"',
)

bash4 = BashOperator(
    task_id="bash4",
    bash_command='echo "Running bash4 - {{ run_id }}"',
)

bash5 = BashOperator(
    task_id="bash5",
    bash_command='echo "Running bash5 - {{ run_id }}"',
)
```

Lunch

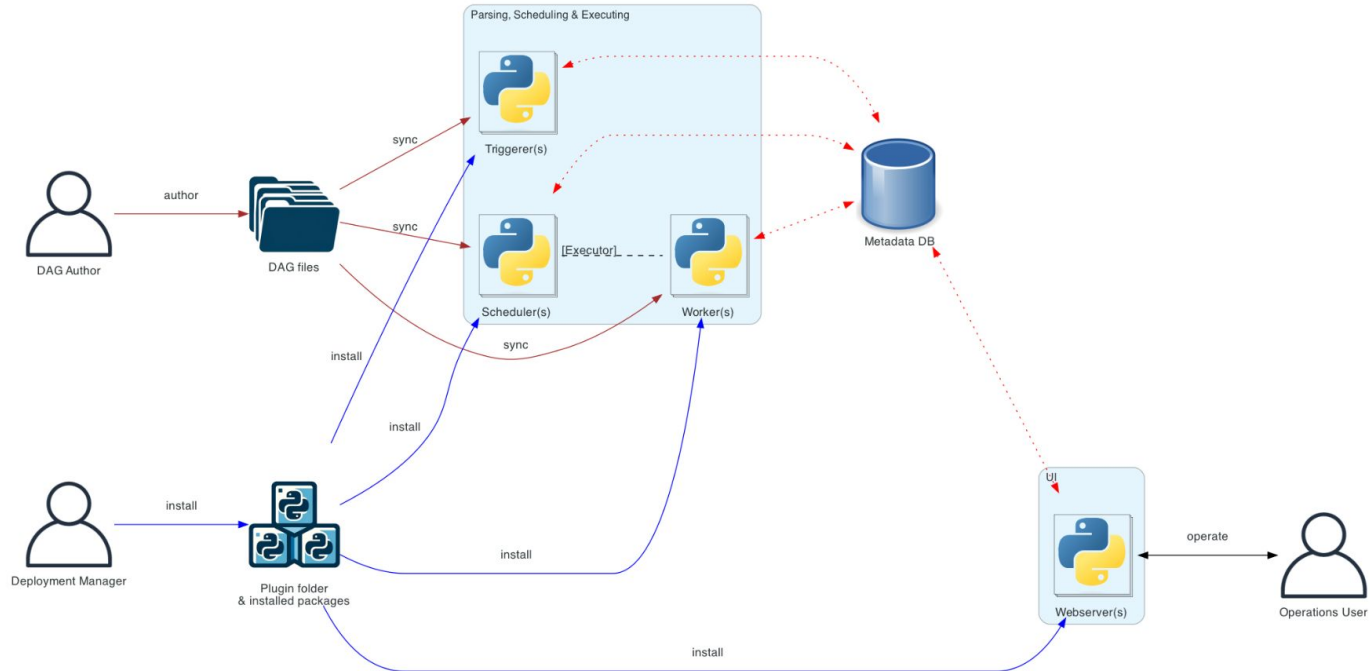
Airflow: Let's dive deeper

- Internal Architecture
- XCom
- Dynamic DAGs
- DAG Dependencies
- DAGs with logic

Airflow Architecture

- **Scheduler**: Triggers scheduled Workflows and submits tasks to the Executor.
- **Executor**: a configuration within the scheduler (not a separate component) determining how tasks are run (e.g., SequentialExecutor, LocalExecutor, CeleryExecutor, etc.).
- **DAG Processor**: Parses DAG files and serializes them into the metadata database.
- **Webserver**: Provides a UI to inspect, trigger, and debug DAGs and tasks.
- **DAG Folder**: Contains DAG files; the scheduler reads this to determine what to run and when.
- **Metadata Database**: Stores the state of workflows and tasks. Required for Airflow to function properly.

Airflow Architecture: Diagram



Execution in Airflow

Scheduler triggers Task instances based on:

- DAG schedule
- Task dependencies
- Execution state (e.g., success/failure)
- Each Task run needs one slot (available Task execution "seat" on a worker)
- Total slots = parallelism limit (per worker and globally)

If no slots are free:

- Tasks are queued until a slot is available
- Sensors can block slots, wasting resources, unless they are deferrable operators

XCOM: Cross-communication

Use case: Task 1 is uploading a file, after which, Task 2 must use that file. Both Tasks are executed asynchronously, but in the dependency order specified. How does Task 1 communicate the file path to Task 2?

XCOM is your answer. It simplifies passing messages between Tasks, but is **specifically for small amounts** of data.

For large amounts of data, cache it and pass that location onward instead.

XCOM: Cross-communication

```
1 @task
2 def get_data() -> str:
3     # Generate sample data
4     df = get_data()
5     filepath = "/tmp/sample_data.csv"
6     df.to_csv(filepath, index=False)
7
8     print(f"Data written to {filepath}")
9     return filepath # XCom value
10
11 @task
12 def load_into_database(csv_path: str):
13     df = pd.read_csv(csv_path)
14     print("Loading into database (simulated):")
15     load_into_postgres(df)
16
17 with DAG(
18     dag_id="xcom_csv_example",
19     start_date=days_ago(1),
20     schedule_interval=None,
21     catchup=False,
22     tags=["example"],
23 ) as dag:
24     filepath = get_data()
25     load_into_database(filepath)
26
```

```
1 def get_data(ti):
2     ...
3     print(f"Data written to {filepath}")
4
5     # Manual XCom push
6     ti.xcom_push(key="csv_path", value=filepath)
7
8 def load_into_database(ti):
9     # Manual XCom pull
10    filepath = ti.xcom_pull(key="csv_path", task_ids="get_data_task")
11    ...
12
13 with DAG(
14     ...
15 ) as dag:
16
17     get_data_task = PythonOperator(
18         task_id="get_data_task",
19         python_callable=get_data,
20     )
21
22     load_task = PythonOperator(
23         task_id="load_into_database",
24         python_callable=load_into_database,
25     )
26
27     get_data_task >> load_task
28
```

Dynamic DAGs: We're writing Python, aren't we?

DAGS generated with Python code are *dynamic*. When used with care, dynamic DAGS can simplify code and add real functionality. When used without care, they add amazing complexity and confusion.

```
1  with DAG("loop_example", ...):
2      first = EmptyOperator(task_id="first")
3      last = EmptyOperator(task_id="last")
4
5      options = ["branch_a", "branch_b", "branch_c", "branch_d"]
6      for option in options:
7          t = EmptyOperator(task_id=option)
8          first >> t >> last
```

DAG dependencies

DAGs are flowcharts but sometimes you want to flowchart your flowcharts.

DAG-to-DAG dependencies are defined by either Triggering or Waiting:

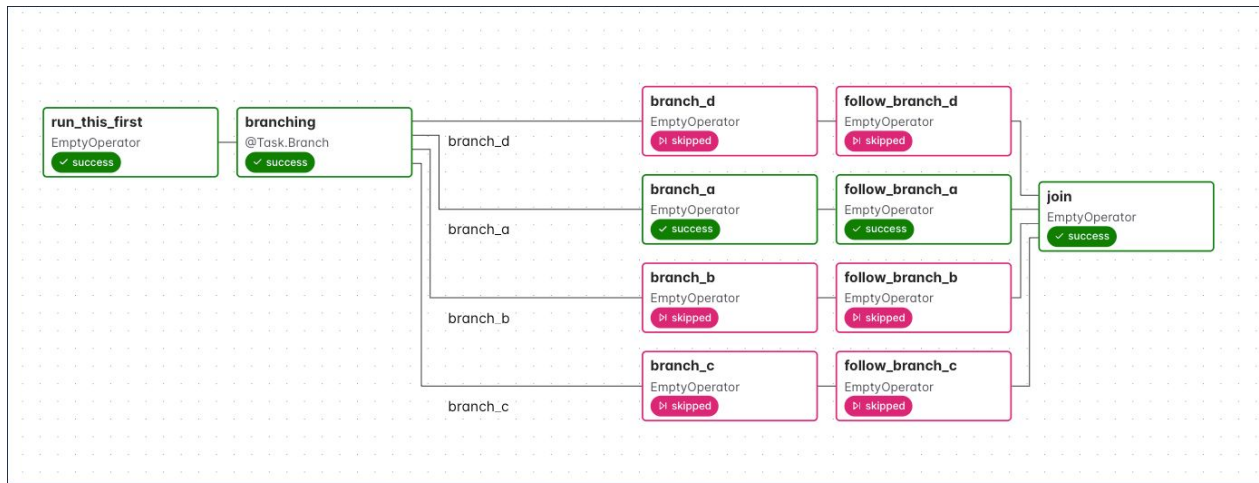
- TriggerDagRunOperator: Start another DAG
- ExternalTaskSensor: Wait for a Task (likely the end of a different DAG)

Warning: This is a loosely-coupled dependency that can turn your nice new tool for transparency and visualization into a nightmare. Use with care and don't nest too deeply.

Logical DAGs: When there are different options

@task.branch (BranchPythonOperator) uses *if* logic to select flowchart Tasks, turning a DAG into a logic tool.

```
1  # from airflow.sdk import task
2
3  result = 1
4
5  @task.branch
6  def choose_branch(result):
7      if result > 0.5:
8          return ['task_a', 'task_b']
9      return ['task_c']
10
11  choose_branch(result)
```



The problem with sensors

- A sensor waits for e.g. 2 hours.
- It occupies a slot for that time.
- It's blocking a slot for that time.

Solution:

Use a Deferrable Operator to free the slot during the wait.

- Deferrable Operators use Triggerer, which waits for events in the background, and, when the condition is met, resumes the Task and uses the slot.
- Triggerer requires an additional lightweight process to be run alongside Airflow, which is why it isn't enabled by default.

Break Time

Custom Operators: Python

Let's write our own Operator using the `@task` decorator. Write a Python function to:

- Import the Iris dataset and store it locally.
- Define the mean of each feature per species
- Show correlation between petal length and petal width
- Define global mean vector (average flower)

```
1  from sklearn.datasets import load_iris
2  import pandas as pd
3
4  # Load as a Bunch object
5  iris = load_iris()
6
7  # Convert to DataFrame
8  df = pd.DataFrame(iris.data, columns=iris.feature_names)
9  df["target"] = iris.target
10
11  print(df.head())
12
```

The **Iris dataset** is one of the most well-known beginner datasets in machine learning. It contains 150 samples of 3 iris species with 4 numeric features (sepal/petal length/width).

Using @task

We can now easily add the @task decorator to turn our Functions into Tasks.

- How would you create the DAG from here?

```
5 @task
6 def import_iris_dataset() -> str:
7     """Download the Iris dataset and save it locally as a CSV file."""
8     iris = load_iris()
9     df = pd.DataFrame(iris.data, columns=iris.feature_names)
10    df["species"] = pd.Categorical.from_codes(iris.target, iris.target_names)
11    filepath = "/tmp/iris.csv"
12    df.to_csv(filepath, index=False)
13    return filepath
14
15 @task
16 def mean_per_species(csv_path: str):
17     """Compute mean of each feature per species."""
18     df = pd.read_csv(csv_path)
19     result = df.groupby("species").mean()
20     print("Mean per species:\n", result)
21
22 @task
23 def petal_length_width_corr(csv_path: str):
24     """Compute correlation between petal length and petal width."""
25     df = pd.read_csv(csv_path)
26     corr = df["petal length (cm)"].corr(df["petal width (cm)"])
27     print(f"Correlation between petal length and width: {corr:.2f}")
28
29 @task
30 def global_mean_vector(csv_path: str):
31     """Compute the global mean of all features."""
32     df = pd.read_csv(csv_path)
33     numeric_cols = df.select_dtypes(include="number")
34     global_mean = numeric_cols.mean()
35     print("Global mean vector:\n", global_mean)
```

Building DAGs: Intuitive Python (but only if you're dutch)

- Why does this work? Does it?

```
1 from airflow import DAG
2 from airflow.utils.dates import days_ago
3
4 ▼ with DAG(
5     dag_id="iris_stats_dag",
6     start_date=days_ago(1),
7     schedule_interval=None,
8     catchup=False,
9 ▼ ) as dag:
10     path = import_iris_dataset()
11     mean_per_species(path)
12     petal_length_width_corr(path)
13     global_mean_vector(path)
14
```

A long time ago in a company far away

Disjointed ETLs: Ingestion and Transformation were handled in 10+ tools

- Ingestion could be using Kubernetes Pods, cronjobs, Pentaho, etc
- Transformation could be in Looker, BigQuery, Scheduled Queries, etc.
- Automations were run from local computers

The plan: Move to Airflow.

Standardisation: What were we doing?

When we took inventory, we were facing the following challenges:

- **Ingestion:** A Kubernetes Pod configured with our own Importer Library.
- **Transformation:** SQL statements that either wrote table-to-table or table-to-view.
 - Needed scheduling and dependency.
- **Custom Pods:** A few special cases, such as A/B testing, machine learning and recommender systems

So...

The Age of Airflow

- We wrote a new Operator/Task: The Importer
 - Configured our importer object and retrieved credentials to authenticate and ingest data
- We wrote a small library: sql-transformation-pods
 - Parsed yaml files with timing and dependencies (grouped in folders) into DAGs and Tasks
- We generalized for other use-cases: A pre-configured GKEPodOperator that allowed us to call for a specific container.

And there we had it: The New Grand Unified Data Platform Tool.

With Airflow, we moved ingestion and transformation from 10 tools to 1, created data visualizations for all pipelines, and centralized credentials and monitoring.

The Age of Airflow

Amusingly, we were using yaml files with SQL statements, dependencies, sometimes unit tests, timing, etc?

Sounds like the **data build tool**

Next week: data build tool (dbt)

They will get angry if you capitalize it. It's in the certification.