



# Apache Airflow

Shani Cohen

# What is Apache Airflow?



# Apache Airflow Use Cases

Automate data ingestion, transformation, and loading for an e-commerce company

**Data Pipelines and ETL**

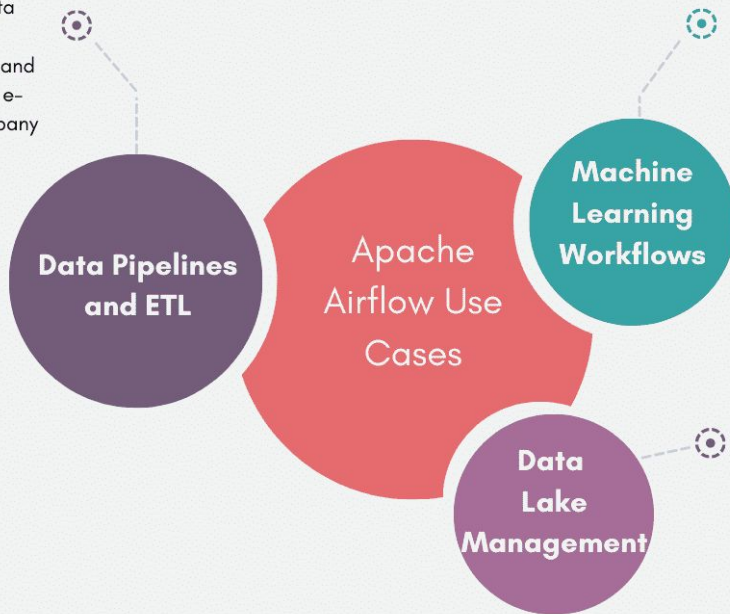
Apache Airflow Use Cases

Manage feature engineering, model training, and deployment for a fraud detection system

**Machine Learning Workflows**

Automate data ingestion, cataloging, and lineage tracking for a large-scale data lake

**Data Lake Management**



# Integrations



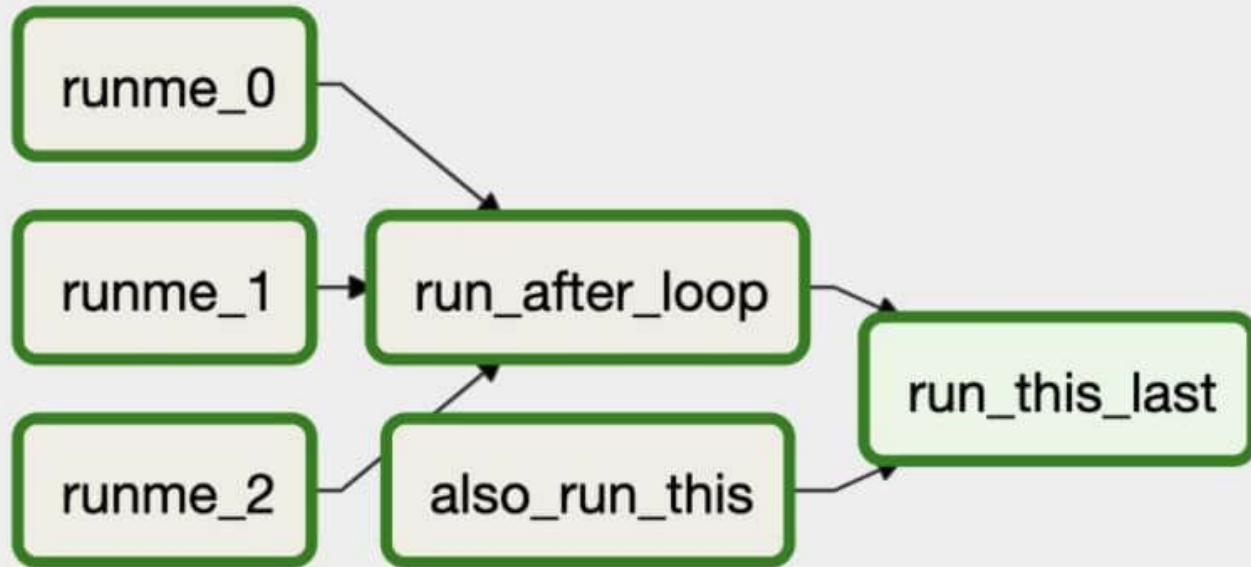
redis



aws

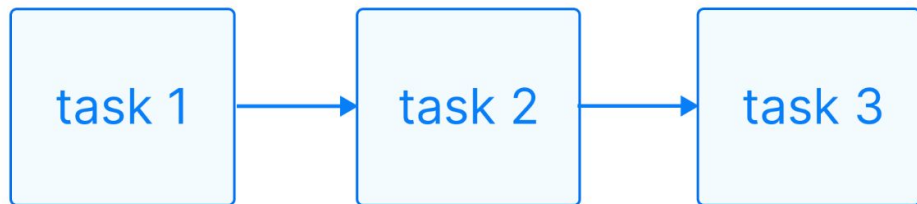


# DAG

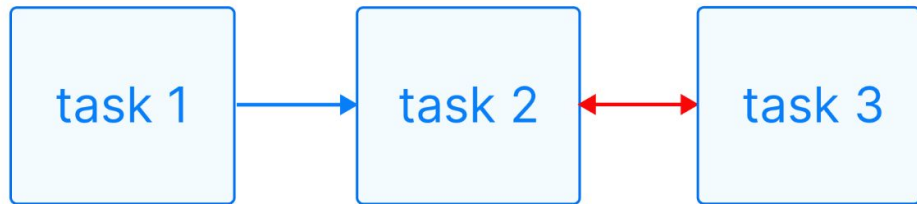


# DAG

Directed

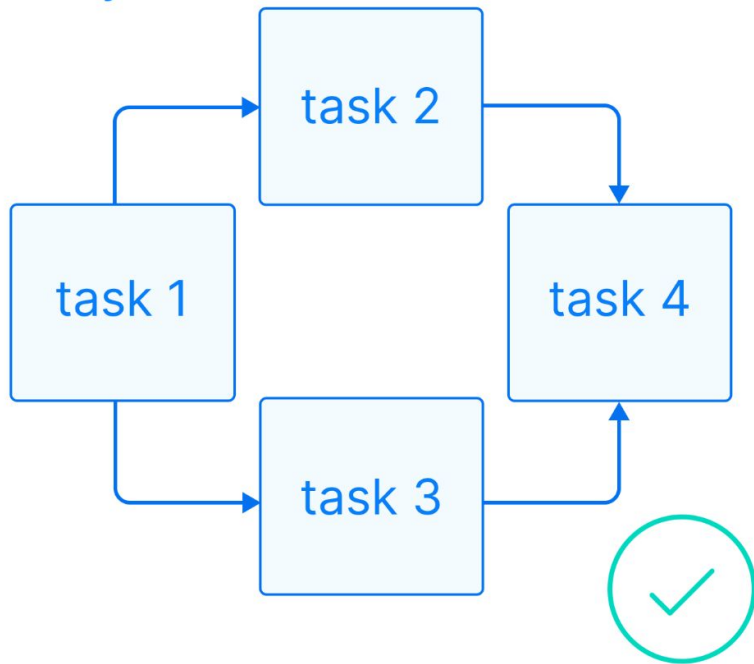


Not directed

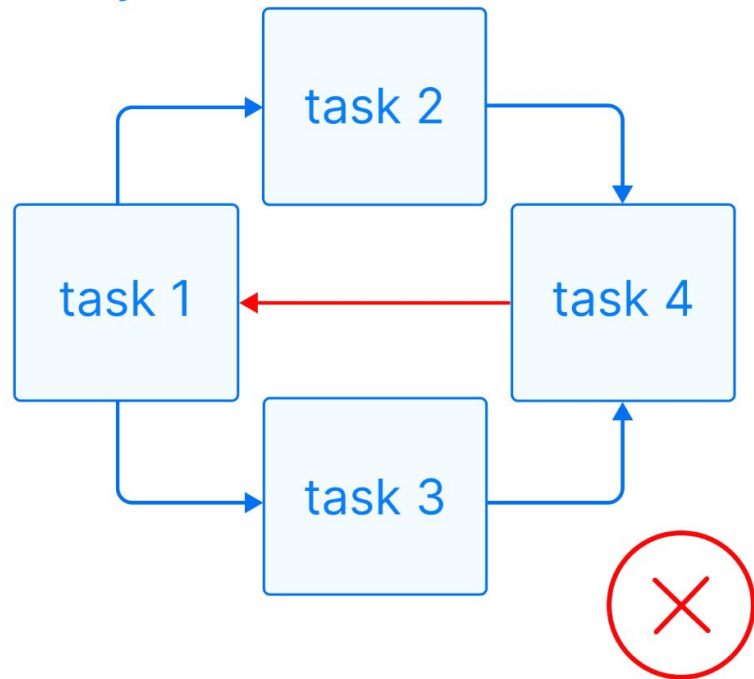


# DAG

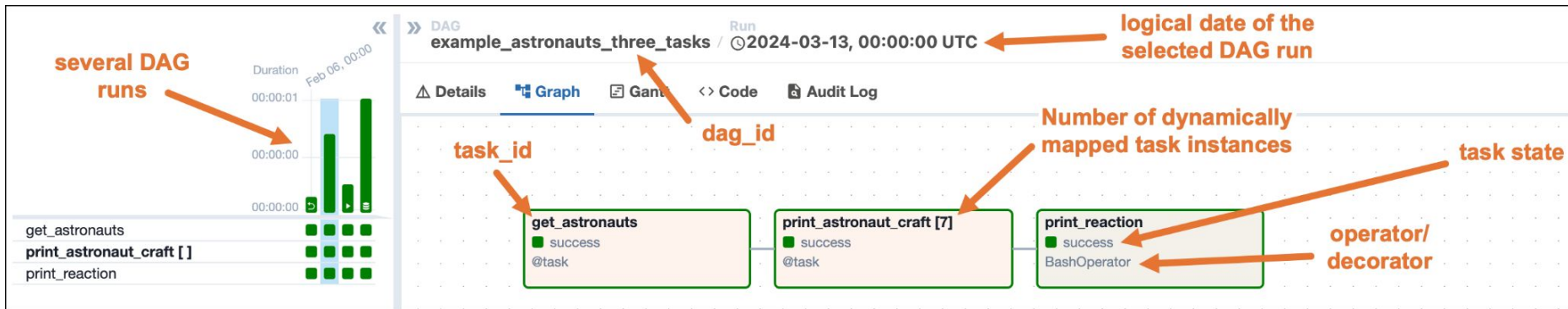
Acyclic



Not Acyclic



# DAG run properties







Grid

Graph

Calendar

Task Duration

Task Tries

Landing Times

Gantt

Details

&lt;&gt; Code

Audit Log

19/09/2024, 06:08:54



25

All Run Types

All Run States

Clear Filters

Auto-refresh

Press **shift** + **/** for Shortcuts

deferred

failed

queued

removed

restarting

running

scheduled

skipped

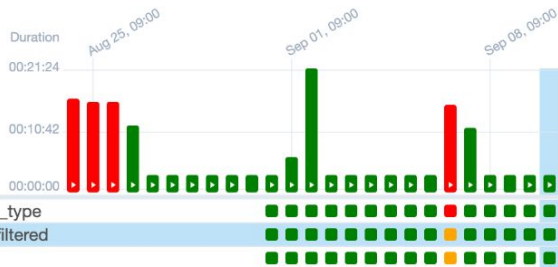
success

up\_for\_reschedule

up\_for\_retry

upstream\_failed

no\_status



DAG

dsc\_model\_to\_type\_frequency /

Run

2024-09-15, 09:00:00 UTC

Task

dsc\_model\_to\_type\_frequency\_filtered

Clear task

Mark state as...

Filter Tasks

Details

Graph

Gantt

&lt;&gt; Code

Logs

XCom

Status success

Task ID dsc\_model\_to\_type\_frequency\_filtered

Run ID manual\_\_2024-09-17T16:58:53+00:00

Operator AthenaOperator

Trigger Rule all\_success

Duration 00:00:31

Started 2024-09-17, 17:00:00 UTC

Ended 2024-09-17, 17:00:31 UTC

# DAG Level Parameters

```
# Define the DAG function & set of parameters
@dag(
    start_date=datetime(2024, 1, 1),
    schedule="@daily",
    catchup=False,
)
def taskflow_dag():
```

- dag\_id
  - start\_date
  - schedule
  - catchup
-

# Workflow as code

```
from datetime import datetime

from airflow import DAG
from airflow.decorators import task
from airflow.operators.bash import BashOperator

# A DAG represents a workflow, a collection of tasks
with DAG(dag_id="demo", start_date=datetime(2022, 1, 1), schedule="0 0 * * *") as dag:
    # Tasks are represented as operators
    hello = BashOperator(task_id="hello", bash_command="echo hello")

    @task()
    def airflow():
        print("airflow")

    # Set dependencies between tasks
    hello >> airflow()
```

# Task

```
t1 = BashOperator(  
    task_id="print_date",  
    bash_command="date",  
)  
  
t2 = BashOperator(  
    task_id="sleep",  
    depends_on_past=False,  
    bash_command="sleep 5",  
    retries=3,  
)
```

# Task Dependencies

```
t1.set_downstream(t2)

# This means that t2 will depend on t1
# running successfully to run.
# It is equivalent to:
t2.set_upstream(t1)

# The bit shift operator can also be
# used to chain operations:
t1 >> t2

# And the upstream dependency with the
# bit shift operator:
t2 << t1

# Chaining multiple dependencies becomes
# concise with the bit shift operator:
t1 >> t2 >> t3

# A list of tasks can also be set as
# dependencies. These operations
# all have the same effect:
t1.set_downstream([t2, t3])
t1 >> [t2, t3]
[t2, t3] << t1
```

# Taskflow

```
from airflow.decorators import task
from airflow.operators.email import EmailOperator

@task
def get_ip():
    return my_ip_service.get_main_ip()

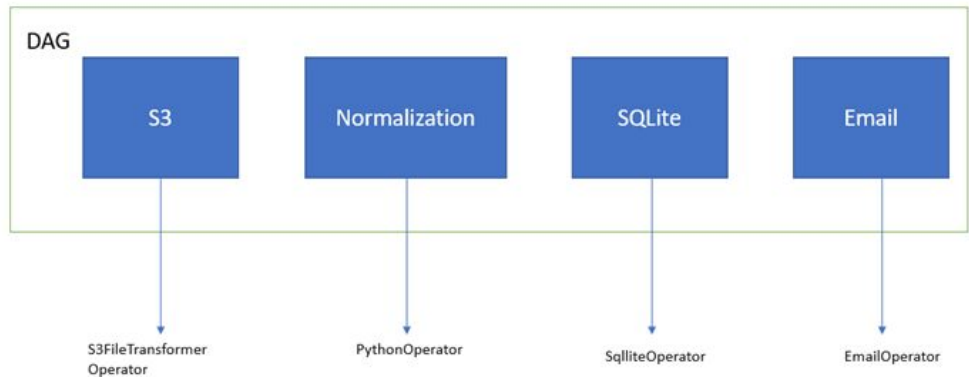
@task(multiple_outputs=True)
def compose_email(external_ip):
    return {
        'subject': f'Server connected from {external_ip}',
        'body': f'Your server executing Airflow is connect

email_info = compose_email(get_ip())
```

# Task Dependencies - Taskflow

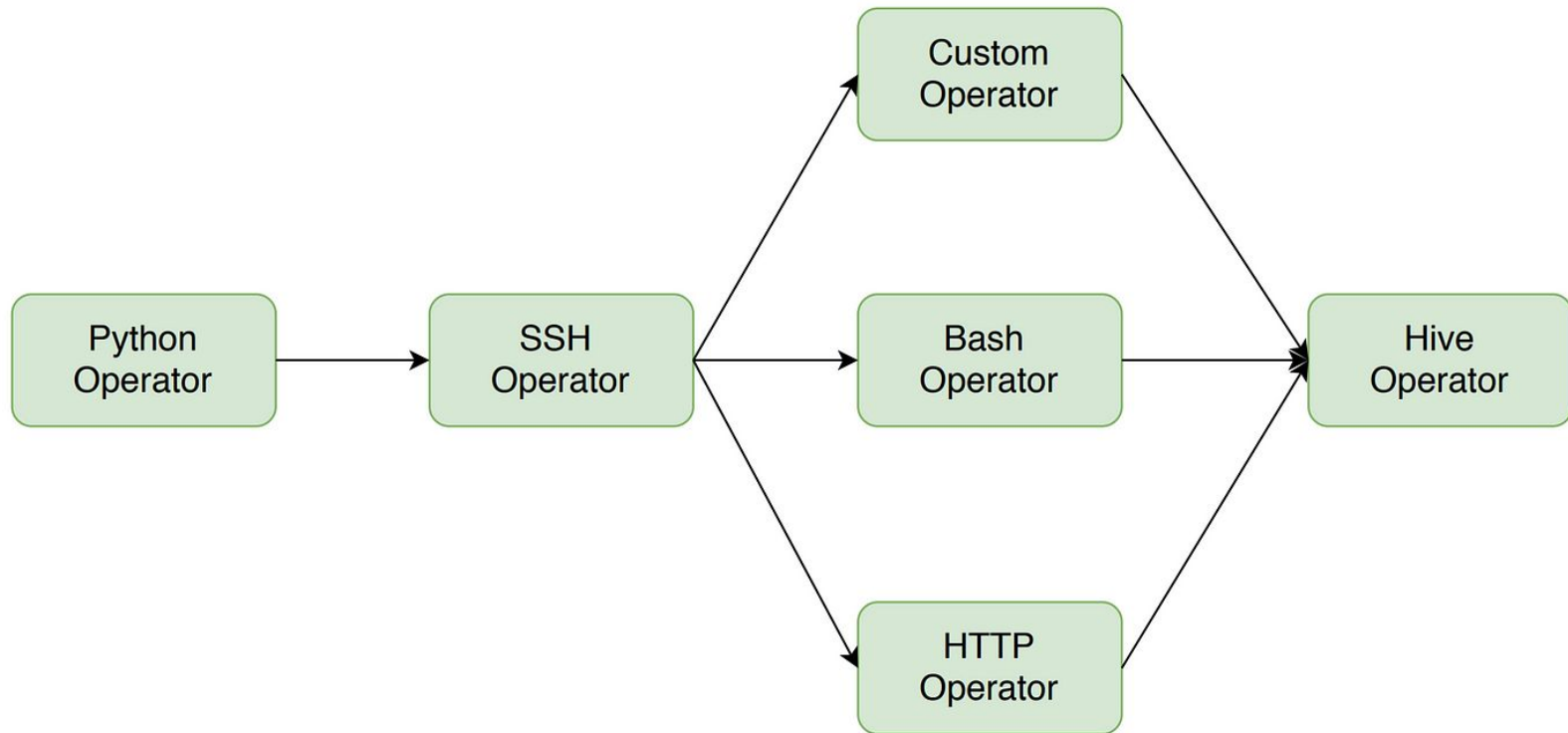
```
order_data = extract()  
order_summary = transform(order_data)  
load(order_summary["total_order_value"])
```

# operators





# Execution Example

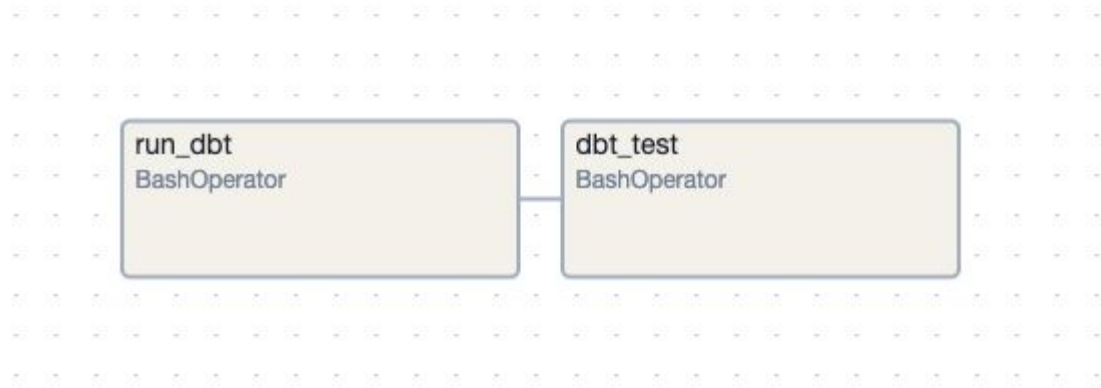


# Workshop Resources:

<https://github.com/shanicohen1902/dbt-cost-allocation-workshop>

## Exercise - Add dbt Test Execution to Your DAG

Modify 'run\_dbt' DAG to add another task that executes  
`dbt test`.



# Jinja

```
templated_command = textwrap.dedent(  
    """  
    {% for i in range(5) %}  
        echo "{{ ds }}"  
        echo "{{ macros.ds_add(ds, 7)}}"  
    {% endfor %}  
    """  
)  
  
t3 = BashOperator(  
    task_id="templated",  
    depends_on_past=False,  
    bash_command=templated_command,  
)
```

# Airflow Templates

```
BashOperator(  
    task_id="print_day_of_week",  
    bash_command="echo Today is {{ execution_date.format('dddd') }}",  
)
```

# Out Of The Box Templats



# Variables

## List Variable

Search ▾

+ Actions ▾ ←

<input type="checkbox"/>	Key ▴ ▾	Val ▴ ▾	Description ▴ ▾	Is Encrypted ▴ ▾
<input type="checkbox"/>	env	dev		False

# Variables

```
from airflow.models import Variable

# Normal call style
foo = Variable.get("foo")

# Auto-deserializes a JSON value
bar = Variable.get("bar", deserialize_json=True)

# Returns the value of default_var (None) if the variable is not set
baz = Variable.get("baz", default_var=None)
```



# Variables

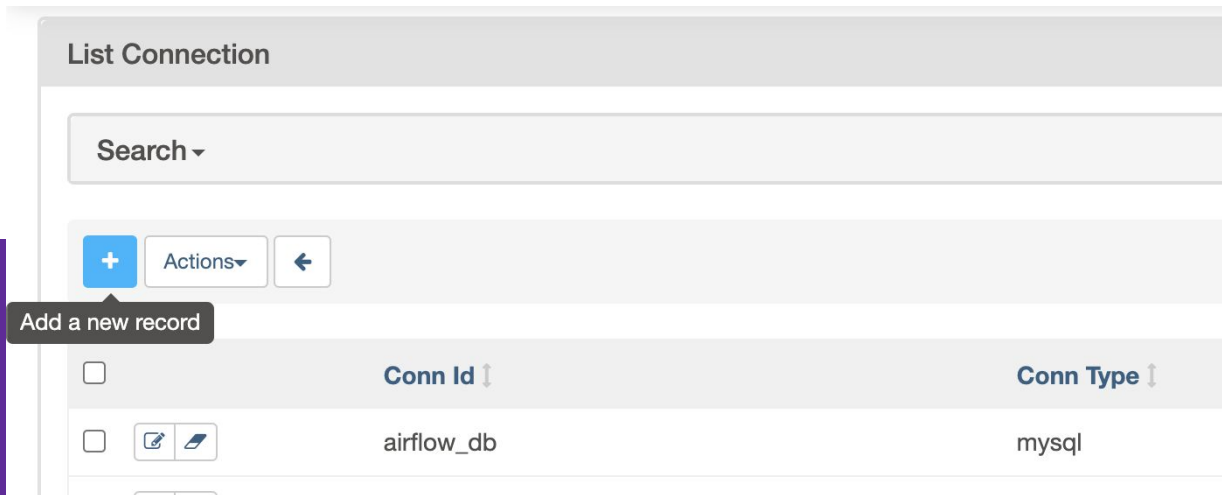
*# Raw value*

```
echo {{ var.value.<variable_name> }}
```

*# Auto-deserialize JSON value*

```
echo {{ var.json.<variable_name> }}
```

# Connections



# Params

## DAG conf Parameters

Names to greet \*:

Linda  
Martha  
Thomas

Define the list of names for which greetings should be generated in the logs. Please have one name per line.

English:



German (Formal):



French:



Generated Configuration JSON and Dagrun Options ▾

# Params

```
from pendulum import datetime
from airflow.decorators import dag, task

@dag(
    start_date=datetime(2023, 6, 1),
    schedule=None,
    catchup=False,
    params={"upstream_color": "Manual run, no upstream color available."},
)
def tdro_example_downstream():
    @task
    def print_color(**context):
        print(context["params"]["upstream_color"])

    print_color()

tdro_example_downstream()
```

# Params

Params are also accessible as a [Jinja template](#) using the `{{ params.my_param }}` syntax.

```
from airflow.utils.template import literal

BashOperator(
    task_id="use_literal_wrapper_to_ignore_jinja_template",
    bash_command=literal("echo {{ params.the_best_number }}"),
)
```

# Params

Params are also accessible as a [Jinja template](#) using the `{{ params.my_param }}` syntax.

```
from airflow.utils.template import literal

BashOperator(
    task_id="use_literal_wrapper_to_ignore_jinja_template",
    bash_command=literal("echo {{ params.the_best_number }}"),
)
```

**Xcom**



# Xcom

```
# pushes data in any_serializable_value into xcom with key "identifier as string"  
task_instance.xcom_push(key="identifier as a string", value=any_serializable_value)
```

```
# pulls the xcom variable with key "identifier as string" that was pushed from within task-1  
task_instance.xcom_pull(key="identifier as string", task_ids="task-1")
```



# Xcom

```
# Pulls the return_value XCOM from "pushing_task"  
value = task_instance.xcom_pull(task_ids='pushing_task')
```

# Xcom

```
SELECT * FROM {{ task_instance.xcom_pull(task_ids='foo', key='table_name') }}
```

Schedule: @once

Next Run ID: 2021-01-01, 00:00:00 UTC



19/09/2024 07:03:09

All Run Types

All Run States

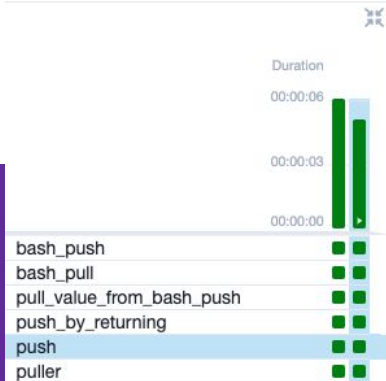
Clear Filters

Auto-refresh

25

Press **shift** + **/** for Shortcuts

deferred
failed
queued
removed
restarting
running
scheduled
shutdown
skipped
success
up\_for\_reschedule
up\_for\_retry
upstream\_failed
no\_status



DAG Run Task  
example\_xcom / 2024-09-19, 07:03:10 UTC / push

Clear task

Mark state as...

Filter DAG by task

Details

Graph

Gantt

Code

Event Log

Logs

XCom

Task Duration

Key

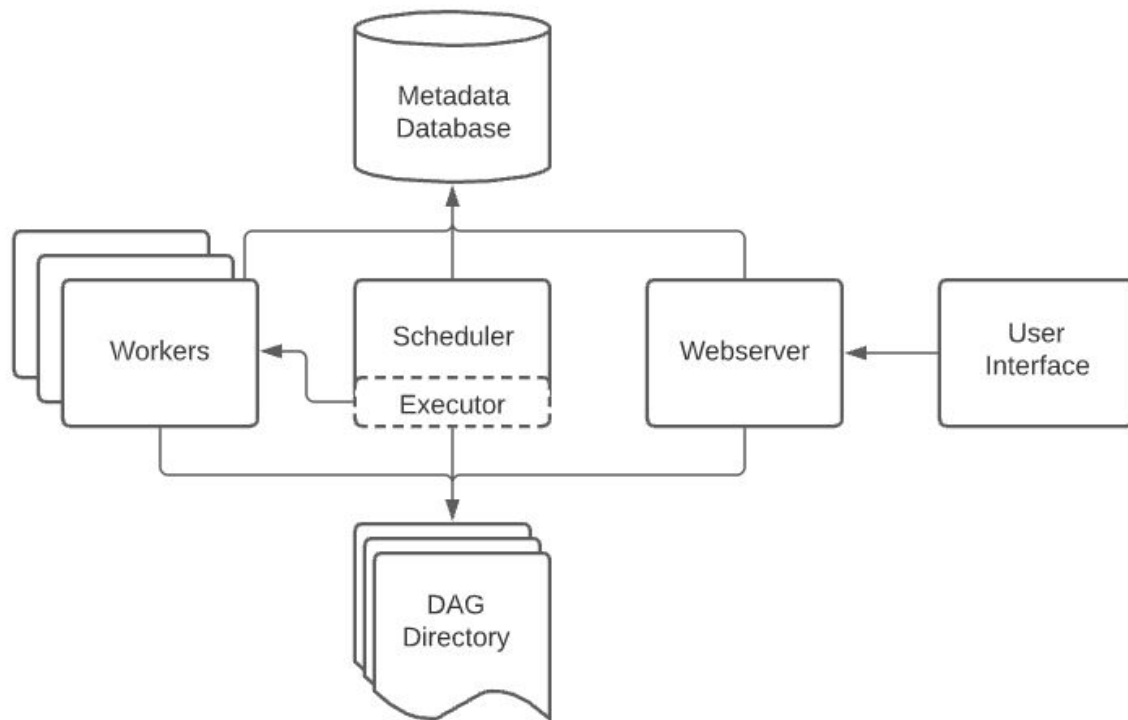
Value

value from pusher 1

[ 3 items  
 0 : 1  
 1 : 2  
 2 : 3  
 1

Copy

# Airflow Architecture



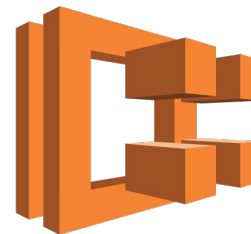
# Executors



C E L E R Y

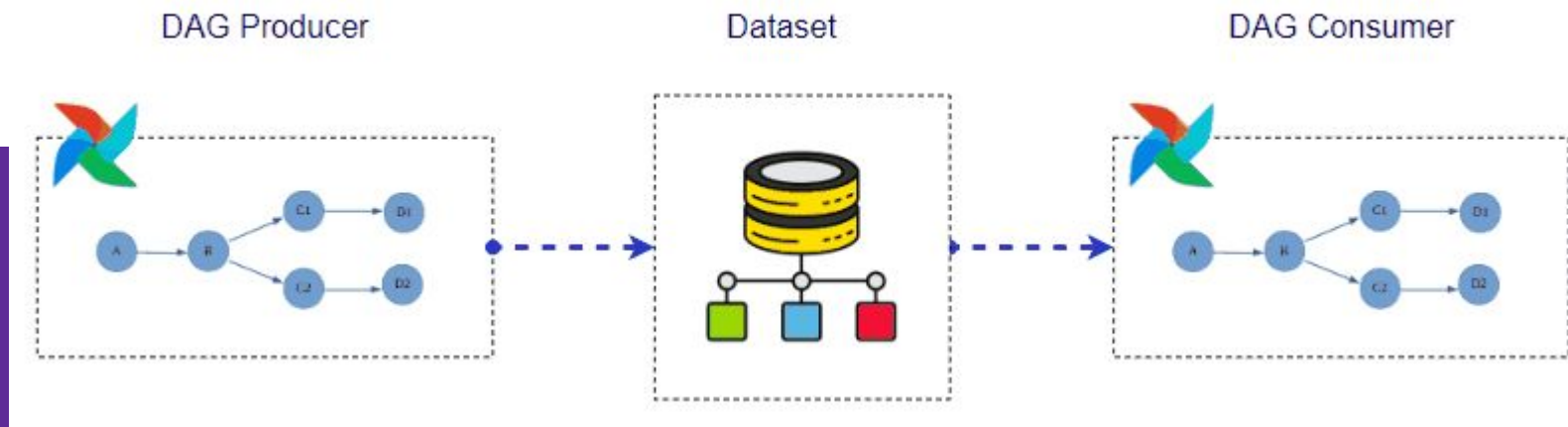


kubernetes



AWS ECS

# Data-aware scheduling



# Data-aware scheduling

```
from airflow.datasets import Dataset

with DAG(...):
    MyOperator(
        # this task updates example.csv
        outlets=[Dataset("s3://dataset-bucket/example.csv")],
        ...,
    )

with DAG(
    # this DAG should be run when example.csv is updated (by dag1)
    schedule=[Dataset("s3://dataset-bucket/example.csv")],
    ...,
):
    ...
```

## Exercise - Xcom + Data-aware scheduling

1. Change 'echo\_airflow\_home' task to return airflow home string to xcom
2. Change 'run\_dbt' task to take airflow\_home sting string from xcom
3. Split run\_dbt into two separate dags - run\_dbt and run\_test. Trigger run\_test every time that run\_dbt successfully executed





**QUIZZZZZZZZ!**

# What does DAG stand for in Apache Airflow?

- A) Data Access Group
- B) Directed Acyclic Graph
- C) Data Aggregation Graph
- D) Dynamic Action Group

---

**Which component is  
responsible for  
executing tasks in  
Airflow?**

- A) Scheduler
- B) Executor
- C) Web Server
- D) Metadata Database

---

# How can tasks communicate with each other in Airflow?

- A) Using global variables
- B) Through direct function calls
- C) By leveraging XComs
- D) Data cannot be passed between tasks

---

# What is the default execution model for Apache Airflow?

- A) LocalExecutor
- B) CeleryExecutor
- C) KubernetesExecutor
- D) SequentialExecutor

---

**Which operator  
would you use to  
execute a Bash  
command in a task?**

- A) PythonOperator
- B) BashOperator
- C) DummyOperator
- D) BranchPythonOperator

---

**What parameter in the DAG definition controls whether past DAG runs should be executed?**

- A) catchup
- B) retries
- C) start\_date
- D) schedule\_interval

---