# Guided Tour to DAG Authoring

Let's flow together

**Workshop**

# Get Airflow Certified

**Thursday, September 21st**
12:00 pm in Trinity 4

Marc Lamberti
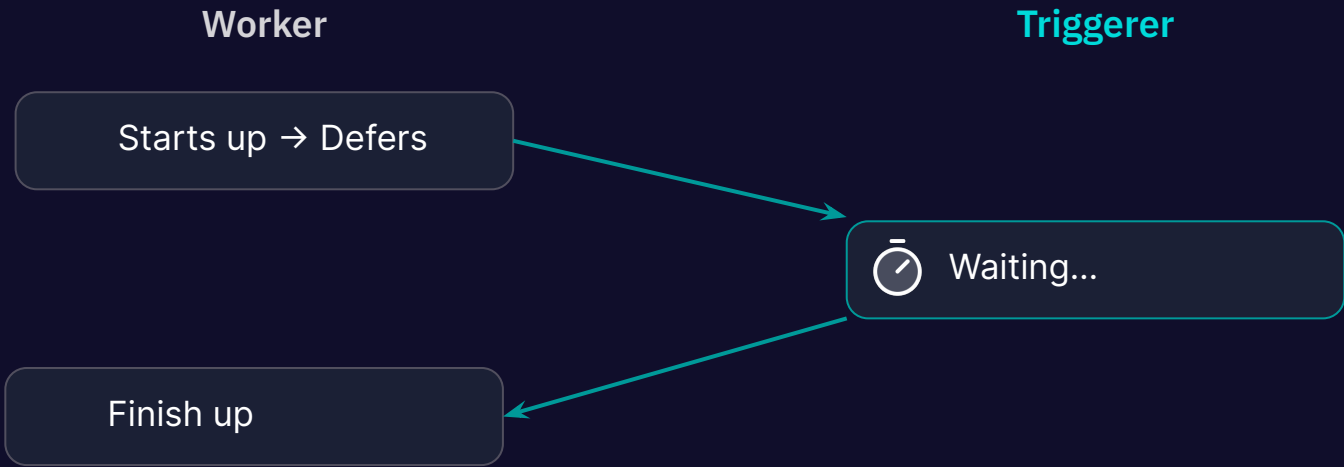Head of Customer Education
at Astronomer

# Agenda

**Deferrable Tasks**
Airflow 2.2

**Dynamic Task Mapping**
Airflow 2.3

**Dynamic DAGs**

**Datasets**
Airflow 2.4

**Setup / Teardown**
Airflow 2.7

**Params**

Deferrable Tasks

# Deferrable Tasks

→ **Airflow 2.2**

→ **Operator/Sensor that can run async**

Don't take up a worker slot

# Deferrable KPO

```python
KubernetesPodOperator(
    task_id="kpo",
    ...,
    deferrable=True,
)
```

# Deferrable KPO

```
$ kubectl get pods -l airflow-worker -w

NAME                STATUS
def-kpo-8czrhyrp    Running
def-kpo-8czrhyrp    Completed

...

def-kpo-69gnklwz    Running
def-kpo-69gnklwz    Completed
```

# Tons of support now!

→ Providers with deferrable support:

- AWS
- Google
- Azure
- DBT
- K8s
- More!

→ `[operators] default_deferrable`

# Custom Operators / Sensors

```
def execute():

    self.defer(
        trigger=SomeTigger(),
        method_name="execute_complete",
    )


def execute_complete():
    return
```

# Custom Operators / Sensors

```
class SomeTrigger(BaseTrigger):

    def serialize(self):

        return ("path.to.SomeTrigger", {})


    async def run(self):

        yield TriggerEvent()
```

# Dynamic Task Mapping

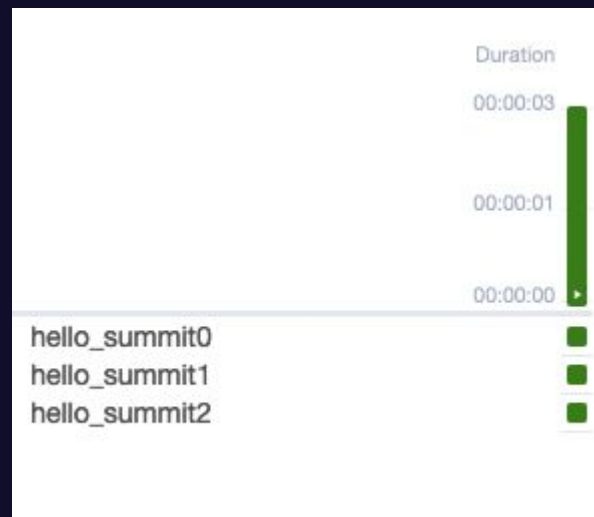→ **Airflow 2.3**
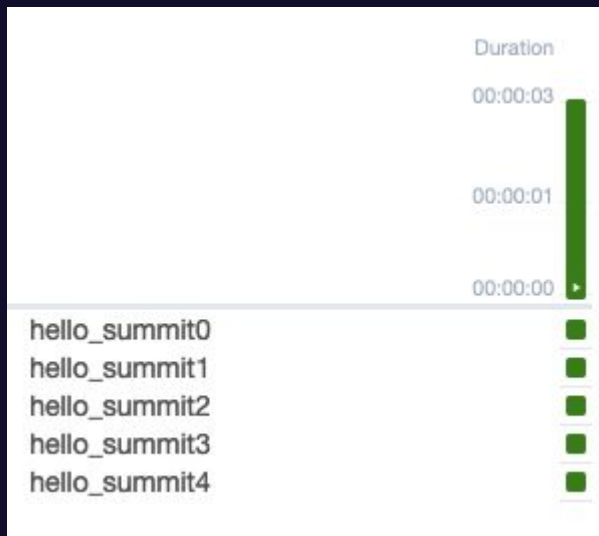
→ **"For loop" for your tasks**

Based on output from a previous task, or static list

→ **"Reduce" tasks**

Task that operates on all results of a mapped task

# Not like this:

```
for file in {s3 bucket}:
        BashOperator({file})
```

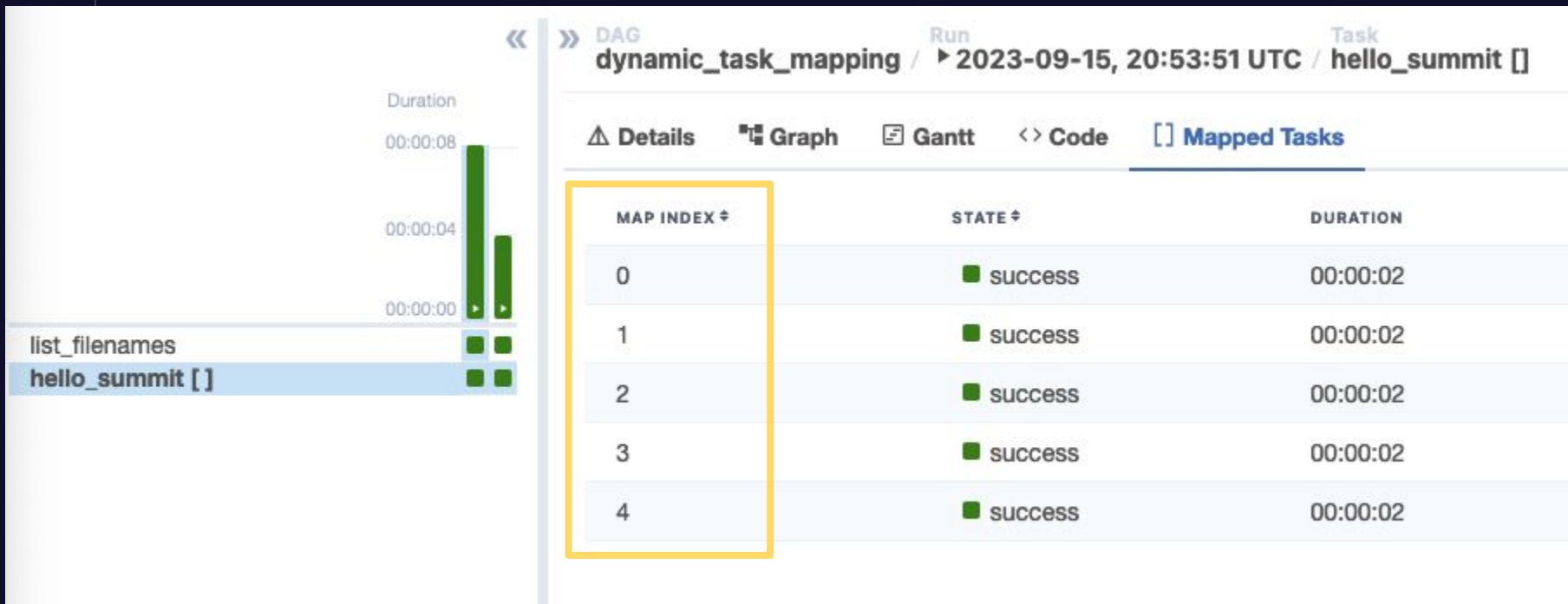# Like this:

```
list_filenames = S3ListOperator(...)


SomeOperator

    .partial(task_id="hello_summit")

    .expand(thing=list_filenames.output)
```
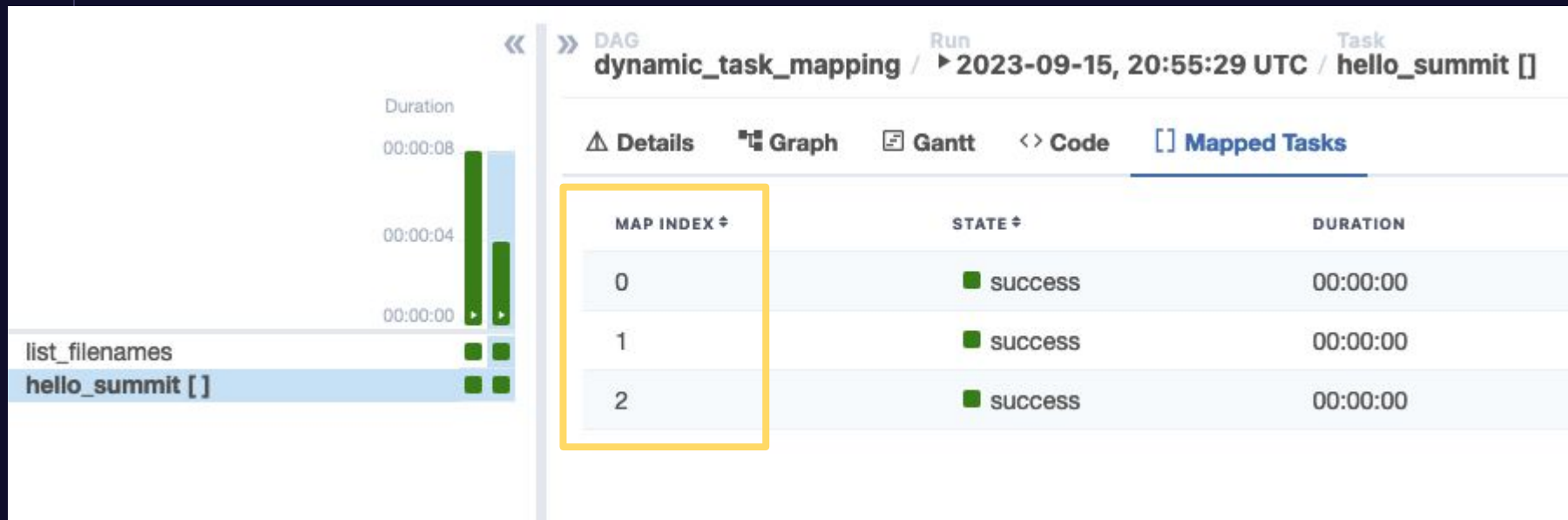
# Like this:

# Like this:

## Reduce:

```
doubled = times_two.expand(x=[1, 2])
sum_(doubled)
```

# Reduce:

# Dynamic DAGs

# Auto Registration

```
for thing in list_of_things:
    with DAG(f"generated_dag_{thing}", ...) as dag:
        …
        globals()[dag_id]
```

# Dynamic DAGs

```
for thing in list_of_things:
    with DAG(f"generated_dag_{thing}", ...):
        ...
```

# Magic Loop in 2.4?

```
desired_id = get_parsing_context().dag_id


for thing in list_of_things:

    dag_id = f"generated_dag_{thing}"
    if desired_id and desired_id != dag_id:
        continue
    ...
```

# Dynamic DAGs

→ Positives:

- Easy code reuse

→ Negatives:

- Debugging complexity
- Scaling

# Datasets

# Datasets

→ **Airflow 2.4**

→ **Data aware scheduling**

Schedule DAG runs based on tasking updating data
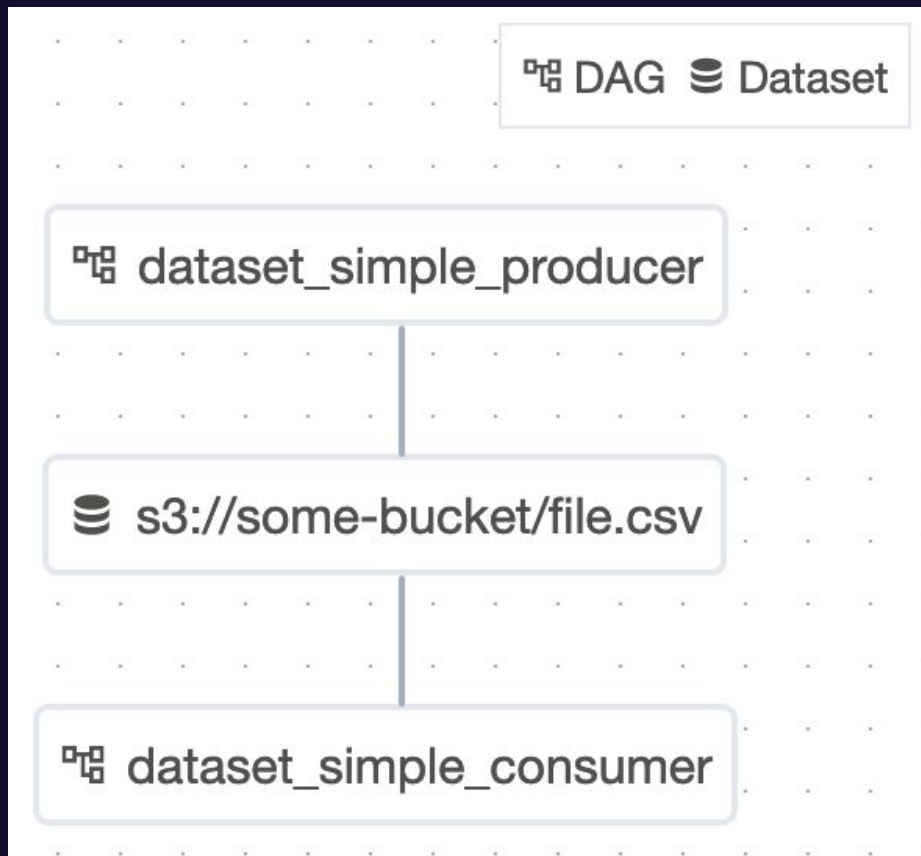
```
MyOperator(
    outlets=[
        Dataset("s3://some-bucket/file.csv")
    ],

    ...,
)
```
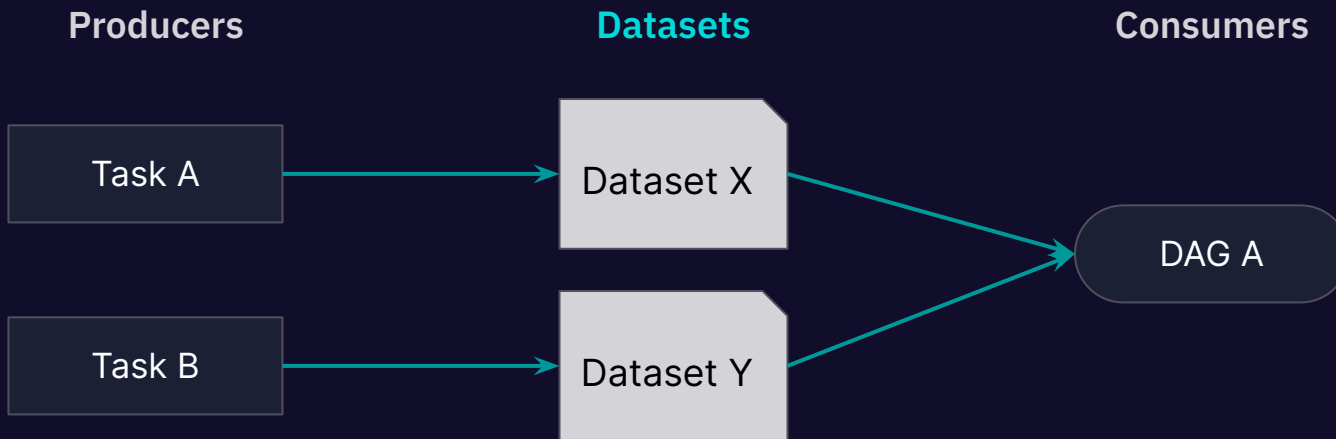
```
with DAG(
    schedule=[
        Dataset("s3://some-bucket/file.csv")
    ],
    ...,
):
    ...
```
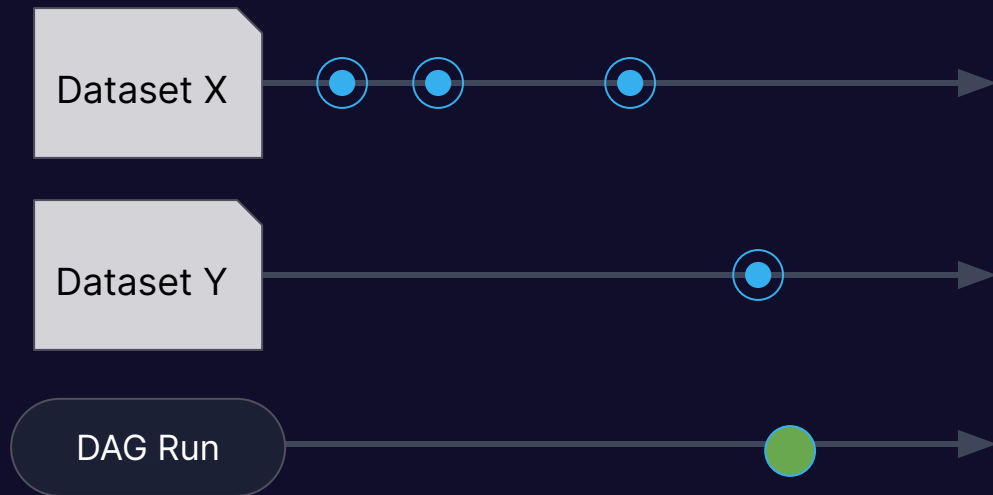
Datasets needed to trigger the next run for

# dataset_simple_consumer2

✕

1 of 2 datasets updated

| Dataset URI | Latest Update |
| --- | --- |
| another-dataset | |
| s3://some-bucket/file.csv | 2023-09-19, 02:49:41 |

Close

→ **What can a Dataset be?**

# Setup / Teardown

# Setup / Teardown

→  **Airflow 2.7**

→  **"Bookend" tasks (support tasks)**

Cleared automatically

Teardown runs if setup ran

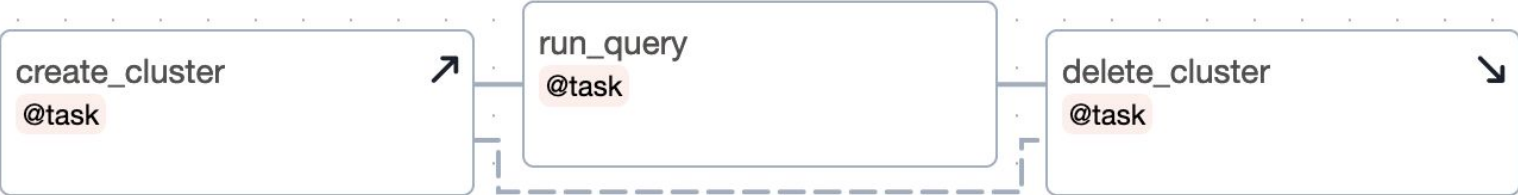Teardown not considered for DAG run state

# Setup / Teardown

```
create_cluster >> run_query >> delete_cluster
```

↓

```
create_cluster >> run_query
run_query >> delete_cluster.as_teardown(setups=create_cluster)
```
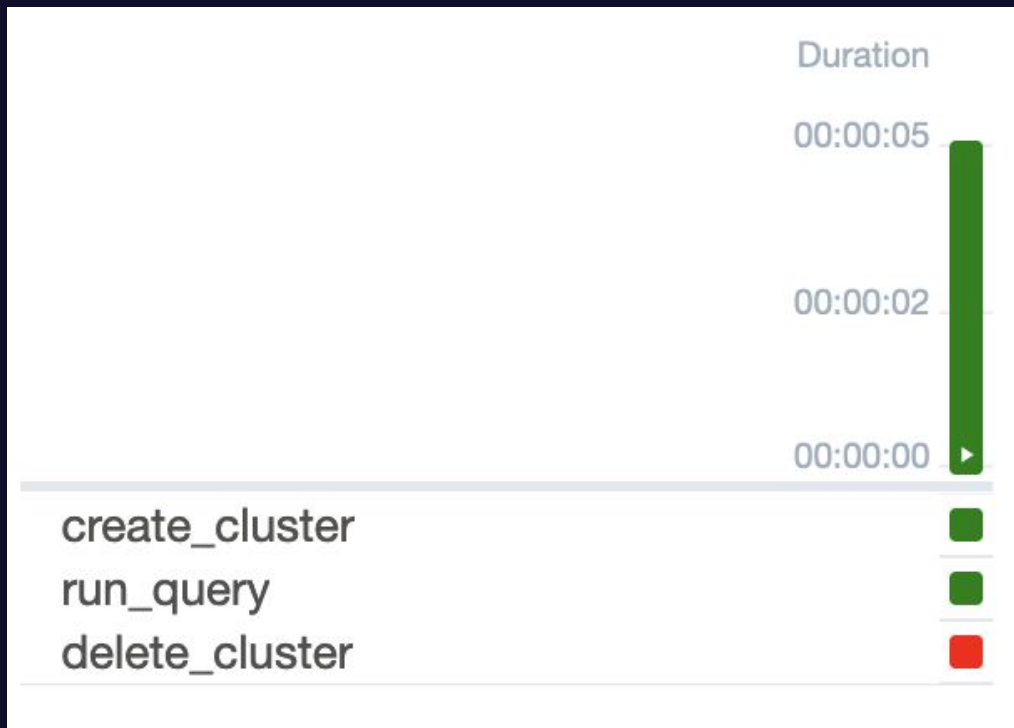
# Setup / Teardown

# Teardown always runs

# DAG Run State

# Clearing

**Clear and Retry** ✕

**Task:** run_query

**Run:** manual__2023-09-18T15:17:08.961246+00:00

Include:

| Past | Future | Upstream | Downstream | Recursive | Only Failed |

**Affected Tasks: 3** ⌃

| TASK NAME ⇅ | MAP INDEX ⇅ | RUN ID ⇅ |
|---|---|---|
| create_cluster | -1 | manual__2023-09-18T15:17:08.961246+00:00 |
| run_query | -1 | manual__2023-09-18T15:17:08.961246+00:00 |
| delete_cluster | -1 | manual__2023-09-18T15:17:08.961246+00:00 |

Cancel                                                              Clear

```
with TaskGroup("my_group") as tg:

    s1 = s1()

    s1 >> w1() >> t1().as_teardown(setups=s1)
```

```
tg >> w2()
```

# Params

# Params

→     **Provide input to runs**

→     JSON Schema

```python
with DAG(

    "some_dag",

    params={

        "rounds": Param(5, type="integer", minimum=3),

    },

):
```

```
def summit(params):

    print(f"Doing {params['rounds']} rounds!")


PythonOperator(task_id="summit", python_callable=summit)
```

```
BashOperator(
        task_id="hello",
        bash_command='echo "Doing $ROUNDS rounds!"',
        env={"ROUNDS": "{{ params.rounds }}"},
    )
```

```
"environment": Param(
    enum=["dev", "stage", "prod"],
    default="stage"
),
```

# Get Involved!

→ Over 2600 contributors

→ All contributions are valuable

→ Join #development-first-pr-support on Slack

Q/A