

Overview

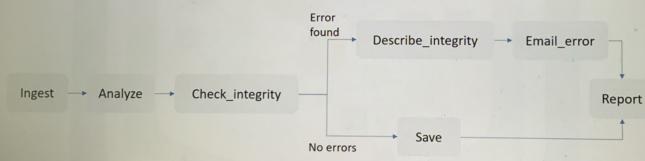
What is Apache Airflow?

- Great open source workflow orchestration tool
- A platform that lets you build and run workflows 
- A workflow is represented as a DAG
- Note: Airflow is not a data streaming solution it is primarily a workflow manager

Unlike Big Data tool such as Apache kafka/storm/Spark

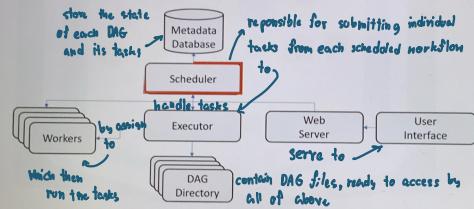
What is Apache Airflow?

DAG specifies the dependencies between tasks, and the order in which execute them
Sample DAG illustrating the labeling of different branches:

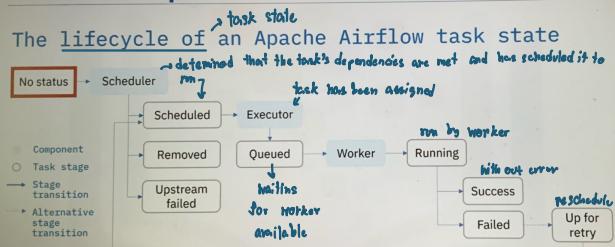


What is Apache Airflow? basic components

Simplified view of Airflow's architecture:



What is Apache Airflow?



Apache Airflow features

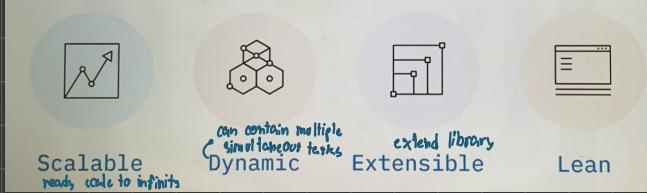


Apache Airflow features



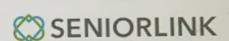
Apache Airflow principles

Airflow pipelines are built on four principles:



Apache Airflow use cases

- Defining and organizing machine learning pipeline dependencies
- Increasing the visibility of batch processes and decoupling them
- Deploying as an enterprise scheduling tool
- Orchestrating SQL transformations in data warehouses
- Sending daily analytics emails



Apache Airflow UI

Objectives

After watching this video, you will be able to:

- Identify current DAGs in your environment
- List different ways to visualize a specific DAG
- Review the code that defines your DAG
- Analyze the duration of each task in your DAG over multiple runs
- Select context metadata for any task instance

DAGs View

This screenshot shows the Airflow DAGs View page. It lists several DAGs: example_bash_operator, example_branch_dag_operator_v3, example_branch_operator, example_complex, example_external_task_marker_child, and example_external_task_marker_parent. Each entry includes the owner (Airflow or airflow), the last run date (e.g., 2021-07-29 00:00:00), the DAG runs count (e.g., 1), and a 'Links' button.

Visualizing your DAG

This screenshot shows the Airflow DAGs View page with a different set of DAGs listed: example_ski_dag, example_subdag_operator, example_trigger_controller_dag, example_xcom, latest_only, latest_only_with_trigger, simple_example, test_utils, and tutorial. The interface is identical to the first screenshot, showing the last run date, DAG runs count, and a 'Links' button.

Visualizing your DAG

This screenshot shows the Airflow DAG view for the 'simple_example' DAG. It displays a Gantt chart with tasks: print_hello, runme_0, runme_1, runme_2, run_after_loop, and run_this_last. The chart shows dependencies between these tasks, such as runme_0 preceding runme_1 and runme_2, and runme_1 preceding run_after_loop.

Visualizing your DAG

This screenshot shows the Airflow DAG view for the 'simple_example' DAG, specifically for a failed run. It displays a Gantt chart with tasks: runme_0, runme_1, runme_2, run_after_loop, and run_this_last. The chart highlights the failure of runme_0 and the subsequent tasks in the loop.

View context metadata for a task

This screenshot shows the Airflow Task Instances view for the 'print_hello' task. It displays the task instance details, including the task ID, run date (2021-07-28T23:01:58.997390+00:00), and the command (print_hello). The log pane shows the output: "Greetings. The date and time are 2021-07-28 23:01:59".

Viewing the code for your DAG

This screenshot shows the Airflow DAG view for the 'simple_example' DAG. It displays the Python code defining the DAG and its tasks:

```
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
import datetime as dt

default_args = {
    'owner': 'me',
    'start_date': dt.datetime(2021, 7, 28),
    'retries': 1,
    'retry_delay': dt.timedelta(minutes=5),
}

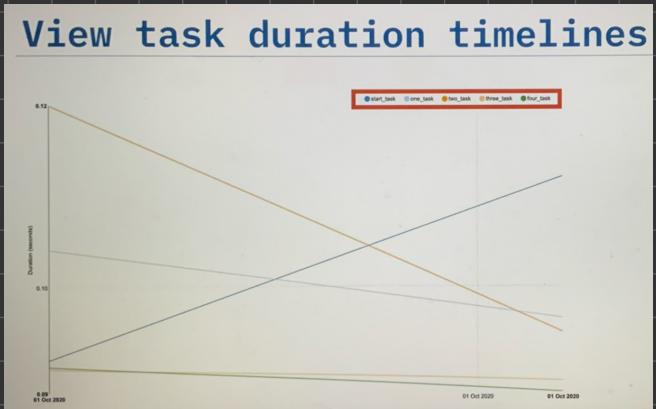
dag = DAG('simple_example',
          description='A simple example DAG',
          default_args=default_args,
          schedule_interval=dt.timedelta(seconds=5),
          )

task1 = BashOperator(
    task_id='print_hello',
    bash_command='echo \'Greetings. The date and time are \'`date`\'\'',
    dag=dag,
)

task2 = BashOperator(
    task_id='print_date',
    bash_command='date',
    dag=dag,
)
task1 >> task2
```

View task duration timelines

This screenshot shows the Airflow Task Duration view for the 'example_bash_operator' DAG. It displays a timeline graph showing the cumulative duration of tasks over time. The x-axis represents dates from 01 Oct 2020 to 02 Oct 2020, and the y-axis represents Duration (seconds).



LAB

Run the command below in the terminal to start Apache Airflow.

```
1 start_airflow
```

Please be patient, it may take upto 5 minutes for airflow to get started.

When airflow starts successfully, you should see an output similar to the one below.

```
Airflow started, waiting for all services to be ready...
Your airflow server is now ready to use and available with username airflow password: MjQ3MzxtbGF2Ym55
You can access your Airflow Webserver at: https://lavanayas-8888.theiadocker-2-labs-prod-theians8s-4-tor81.pro
[2024-05-01] [INFO] [airflow.main:100] -
```

- List DAGs: airflow dags list
- List Tasks: airflow tasks list example_bash_operator
- Run a specific task: airflow tasks test example_bash_operator runme_1 2024-05-01

Exercise 3 - List all DAGs

Apache airflow gives us some handy command line options to work with.

Run the command below in the terminal to list out all the existing DAGs.

```
1 airflow dags list
```

Exercise 4 - List tasks in a DAG

Run the command below in the terminal to list out all the tasks in the DAG named `example_bash_operator`.

```
1 airflow tasks list example_bash_operator
```

Here `example_bash_operator` is the name of the DAG.

Try listing out the tasks for the DAG `tutorial`.

Exercise 5 - Pause/Unpause a DAG

Run the command below in the terminal to unpause a DAG.

```
1 airflow dags unpause tutorial
```

Here `tutorial` is the name of the DAG.

Pause a DAG.

```
1 airflow dags pause tutorial
```

Try to unpause the DAG named `example_branch_operator`.

Build DAG using Airflow

Objectives

After watching this video, you will be able to:

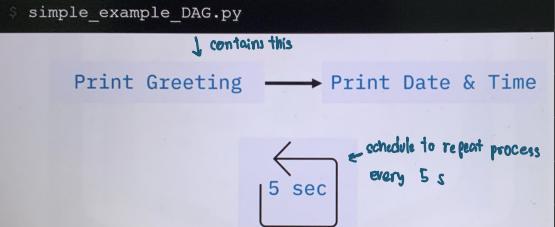
- Interpret an Airflow pipeline as a Python script that defines an Airflow DAG object
- List the key components of a DAG definition file
- Create tasks by instantiating operators in your DAG definition file
- Set up dependencies amongst tasks

Airflow DAG definition script

Python script blocks:

- Library imports
- DAG arguments
- DAG definition
- Task definitions
- Task pipeline

Build an Airflow pipeline



Airflow pipeline script

Python script blocks:

- Python library imports

```
from airflow import DAG for Bash
from airflow.operators.bash_operator import BashOperator
import datetime as dt → for time to schedule
```

Airflow pipeline script

Python script containing:

- Python library imports
- DAG arguments

```
default_args = {
    'owner': 'me',
    'start_date': dt.datetime(2021, 7, 28),
    'retries': 1, → number of times it should keep trying
    'retry_delay': dt.timedelta(minutes=5), → time to wait between subsequent tries → 5 min
}
```

if it failing
↓
This ex is one
if it does fail

Airflow pipeline script

Python script containing:

- Python library imports
- DAG arguments
- DAG definition

```
dag=DAG('simple_example', → name of your DAG
        description='A simple example DAG', →
        default_args=default_args, → default arguments to apply
        schedule_interval=dt.timedelta(seconds=5), →
        ) → run repeatedly
```

Airflow pipeline script

Task Definitions

```
task1 = BashOperator(
    task_id='print_hello', → respective ID
    bash_command='echo \'Greetings. The date and time are \'', → call bash
    dag=dag, → assign DAG Block
)
task2 = BashOperator(
    task_id='print_date',
    bash_command='date',
    dag=dag,
)
```

Airflow pipeline script

Python script containing:

- Python library imports
- DAG arguments
- DAG definition
- Task definitions
- Task pipeline

task1 >> task2
↓ mean
'print_hello' 'print_date'
↓ This one runs first then

LAB

Start Apache Airflow in the lab environment.

```
1 start_airflow
```

Please be patient, it will take a few minutes for airflow to get started.

Exercise 3 - Explore the anatomy of a DAG

An Apache Airflow DAG is a python program. It consists of these logical blocks.

- Imports
- DAG Arguments
- DAG Definition
- Task Definitions
- Task Pipeline

A typical `imports` block looks like this.

```
1 # import the libraries
2
3 from datetime import timedelta
4 # The DAG object; we'll need this to instantiate a DAG
5 from airflow.models import DAG
6 # Operators; you need this to write tasks!
7 from airflow.operators.bash_operator import BashOperator
8 # This makes scheduling easy
9 from airflow.utils.dates import days_ago
```

A typical `DAG Arguments` block looks like this.

```
1 #defining DAG arguments
2
3 # You can override them on a per-task basis during operator
4 default_args = {
5     'owner': 'Lavanya',
6     'start_date': days_ago(0),
7     'email': ['lavanya@someemail.com'],
8     'retries': 1,
9     'retry_delay': timedelta(minutes=5),
10 }
```

to retry again

DAG arguments are like settings for the DAG.

The above settings mention

- the owner `name`,
- when this DAG should run from: `days_ago(0)` means today,
- the email address where the alerts are sent to,
- the number of retries in case of failure, and
- the time delay between retries.

The other options that you can include are:

- 'queue' : The name of the queue the task should be a part of
- 'pool' : The pool that this task should use.
- 'email_on_failure' : Whether an email should be sent to the owner on failure
- 'email_on_retry' : Whether an email should be sent to the owner on retry
- 'priority_weight' : Priority weight of this task against other tasks.
- 'end_date' : End date for the task
- 'wait_for_downstream' : Boolean value indicating whether it should wait for downtime
- 'sla' : Time by which the task should have succeeded. This can be a timedelta object.
- 'execution_timeout' : Time limit for running the task. This can be a timedelta object.
- 'on_failure_callback' : Some function, or list of functions to call on failure
- 'on_success_callback' : Some function, or list of functions to call on success
- 'on_retry_callback' : Another function, or list of functions to call on retry
- 'sla_miss_callback' : Yet another function, or list of functions when sla is missed
- 'on_skipped_callback' : Some function to call when the task is skipped
- 'trigger_rule' : Defines the rule by which the generated task gets triggered

A typical `DAG definition` block looks like this.

```
1 # define the DAG
2 dag = DAG(
3     dag_id='sample-etl-dag',
4     default_args=default_args,
5     description='Sample ETL DAG using Bash',
6     schedule_interval=timedelta(days=1),
7 )
```

Here you are creating a variable named dag by instantiating the `DAG` class with the following parameters.

`sample-etl-dag` is the ID of the DAG. This is what you see on the web console.

you are passing the dictionary `default_args`, in which all the defaults are defined.

`description` helps us in understanding what this DAG does.

`schedule_interval` tells us how frequently this DAG runs. In this case every day. (`days=1`).

A typical `task definitions` block looks like this:

```
1 # define the tasks
2
3 # define the first task named extract
4 extract = BashOperator(
5     task_id='extract',
6     bash_command='echo "extract"',
7     dag=dag,
8 )
9
10 # define the second task named transform
11 transform = BashOperator(
12     task_id='transform',
13     bash_command='echo "transform"',
14     dag=dag,
15 )
16
17 # define the third task named load
18
19 load = BashOperator(
20     task_id='load',
21     bash_command='echo "load"',
22     dag=dag,
23 )
```

A task is defined using:

- A `task_id` which is a string and helps in identifying the task.
- What bash command it represents.
- Which dag this task belongs to.

A typical `task pipeline` block looks like this:

```
1 # task pipeline
2 extract >> transform >> load
```

Task pipeline helps us to organize the order of tasks.

Here the task `extract` must run first, followed by `transform`, followed by the task `load`.

Exercise 4 - Create a DAG

Let us create a DAG that runs daily, and extracts user information from `/etc/passwd` file, transforms it, and loads it into a file.

This DAG has two tasks `extract` that extracts fields from `/etc/passwd` file and `transform_and_load` that transforms and loads data into a file.

```
1 # import the libraries
2
3 from datetime import timedelta
4 # The DAG object; we'll need this to instantiate a DAG
5 from airflow.models import DAG
6 # Operators; you need this to write tasks!
7 from airflow.operators.bash_operator import BashOperator
8 # This makes scheduling easy
9 from airflow.utils.dates import days_ago
10
11 #defining DAG arguments
12
13 # You can override them on a per-task basis during operat
14 default_args = {
15     'owner': 'your_name_here',
16     'start_date': days_ago(0),
17     'email': ['your_email_here'],
18     'retries': 1,
19     'retry_delay': timedelta(minutes=5),
20 }
21
22 # defining the DAG
23
24 # define the DAG
25 dag = DAG(
26     'my-first-dag',
27     default_args=default_args,
28     description='My first DAG',
29     schedule_interval=timedelta(days=1),
30 )
31
32 # define the tasks
33
34 # define the first task
35
36 extract = BashOperator(
37     task_id='extract',
38     bash_command='cut -d ":" -f1,3,6 /etc/passwd > /home/passwd > /home/project/airflow/dags/extracted-data.txt',
39     dag=dag,
40 )
41
42 # define the second task
43 transform_and_load = BashOperator(
44     task_id='transform',
45     bash_command='tr ":" "," < /home/project/airflow/dags/extracted-data.txt > /home/project/airflow/dags/transformed-data.txt',
46     dag=dag,
47 )
48
49 # task pipeline
50 extract >> transform_and_load
```

Create a new file by choosing File->New File and name it `my_first_dag.py`. Copy the code above and paste it into `my_first_dag.py`.

Exercise 5 - Submit a DAG

Submitting a DAG is as simple as copying the DAG python file into `dags` folder in the `AIRFLOW_HOME` directory.

Airflow searches for Python source files within the specified `DAGS_FOLDER`. The location of `DAGS_FOLDER` can be located in the `airflow.cfg` file, where it has been configured as `/home/project/airflow/dags`.

```
airflow > airflow.cfg
1 [core]
2 # The folder where your air-flow pipelines live, most likely a
3 # subfolder in a code repository. This path must be absolute.
4 dags_folder = /home/project/airflow/dags
```

Airflow will load the Python source files from this designated location. It will process each file, execute its contents, and subsequently load any DAG objects present in the file.

Therefore, when submitting a `DAG`, it is essential to position it within this directory structure. Alternatively, the `AIRFLOW_HOME` directory, representing the structure `/home/project/airflow`, can also be utilized for DAG submission.

```
theia@theiadocker-lavanya:~/home/project >
theia@theiadocker-lavanya:~/home/project$ echo $AIRFLOW_HOME
/home/project/airflow
```

Open a terminal and run the command below to submit the DAG that was created in the previous exercise.

```
1 cp my_first_dag.py $AIRFLOW_HOME/dags
```

Verify that your DAG actually got submitted.

Run the command below to list out all the existing DAGs.

```
1 airflow dags list
```

Verify that `my-first-dag` is a part of the output.

```
1 airflow dags list|grep "my-first-dag"
```

You should see your DAG name in the output.

Run the command below to list out all the tasks in `my-first-dag`.

```
1 airflow tasks list my-first-dag
```

You should see 2 tasks in the output.

LAB

Exercise 2 - Create a basic shell script

Here we will define a shell script `extract_transform_load.sh` which will define a pipeline of tasks such as

- extract
- transform
- load

For now, let the shell script have basic echo statements for `extract`, `transform` and `load`.

```
1 #!/bin/bash
2
3 echo "Extract data"
4
5 echo "Transform data"
6
7 echo "Load data"
```

In the next section we will define a `DAG` to call and execute the shell script.

A typical `task definitions` block looks like this:

```
1 # define the tasks
2
3 # define the task named extract_transform_and_load to call the shell script
4 extract_transform_and_load = BashOperator(
5     task_id="extract_transform_and_load",
6     bash_command='/home/project/airflow/dags/extract_transform_load.sh',
7     dag=dag,
8 )
```

A task is defined using:

- A `task_id` which is a string and helps in identifying the task.
- What bash command it represents. Here we are calling the shell script `extract_transform_load.sh` which we previously defined.
- Which dag this task belongs to.

A typical `task pipeline` block looks like this:

```
1 # task pipeline
2 extract_transform_and_load
```

When we execute the task `extract_transform_and_load` the code in the shell script gets executed.

Exercise 3 - Explore the anatomy of a DAG

An Apache Airflow DAG is a python program. It consists of these logical blocks.

- Imports
- DAG Arguments
- DAG Definition
- Task Definitions
- Task Pipeline

A typical `imports` block looks like this.

```
1 # import the libraries
2
3 from datetime import timedelta
4 # The DAG object; we'll need this to instantiate a DAG
5 from airflow import DAG
6 # Operators; we need this to write tasks!
7 from airflow.operators.bash_operator import BashOperator
8 # This makes scheduling easy
9 from airflow.utils.dates import days_ago
```

A typical `DAG Arguments` block looks like this.

```
1 #defining DAG arguments
2
3 # You can override them on a per-task basis during operator initialization
4 default_args = {
5     'owner': 'Ramesh Sannareddy',
6     'start_date': days_ago(0),
7     'email': ['ramesh@somedomain.com'],
8     'email_on_failure': True,
9     'email_on_retry': True,
10    'retries': 1,
11    'retry_delay': timedelta(minutes=5),
12 }
```

DAG arguments are like settings for the DAG.

The above settings mention

- the `owner` name,
- when this DAG should run from: `days_ago(0)` means today,
- the email address where the alerts are sent to,
- whether alert must be sent on failure,
- whether alert must be sent on retry,
- the number of retries in case of failure, and
- the time delay between retries.

Exercise 4 - ETL process on a `/etc/passwd` file

Here we will first

- Create a new shell script called `my_first_dag.sh` to perform the ETL process.
- Create a `DAG` file `my_first_dag.py` which will run daily and defines a task `etl` to call the shell script `my_first_etl.sh`.
- Submit the `DAG`

Create a new shell script `my_first_dag.sh` by selecting `File->New File`.

- The shell script extracts the first, third and sixth fields from `/etc/passwd` file using the `cut` command and writes to a new file `extracted-data.txt`.
- Next the `extracted-data.txt` is transformed to a `.csv` file and loaded into a new file called `transformed-data.csv` using `tr` command.

Copy the code below in the shell script.

```
1 #!/bin/bash
2 echo "extract_transform_and_load"
3 cut -d ":" -f1,3,6 /etc/passwd > /home/project/airflow/dags/extracted-data.txt
4
5 tr ";" "," < /home/project/airflow/dags/extracted-data.txt > /home/project/airflow/dags/transformed-data.csv
```

Click `Save` to save the shell script.

Create a new `DAG` file `my_first_dag.py` by selecting `File->New File`.

This DAG has one task `etl` that calls the shell script `my_first_dag.sh`

```
1 # import the libraries
2
3 from datetime import timedelta
4 # The DAG object; we'll need this to instantiate a DAG
5 from airflow import DAG
6 # Operators; we need this to write tasks!
7 from airflow.operators.bash_operator import BashOperator
8 # This makes scheduling easy
9 from airflow.utils.dates import days_ago
10
11 #defining DAG arguments
12
13 # You can override them on a per-task basis during operator initialization
14 default_args = {
15     'owner': 'Ramesh Sannareddy',
16     'start_date': days_ago(0),
17     'email': ['ramesh@somedomain.com'],
18     'email_on_failure': False,
19     'email_on_retry': False,
20     'retries': 1,
21     'retry_delay': timedelta(minutes=5),
22 }
```

```
24 # defining the DAG
25
26 # define the DAG
27 dag = DAG(
28     'my-first-dag',
29     default_args=default_args,
30     description='My first DAG',
31     schedule_interval=timedelta(days=1),
32 )
33
34 # define the task **extract_transform_and_load** to call shell script
35
36 #calling the shell script
37 extract_transform_load = BashOperator(
38     task_id="extract_transform_load",
39     bash_command="/home/project/airflow/dags/my_first_dag.sh",
40     dag=dag,
41 )
42
43 # task pipeline
44 extract_transform_load
```

A typical `DAG definition` block looks like this.

```
1 # define the DAG
2 dag = DAG(
3     dag_id='sample-etl-dag',
4     default_args=default_args,
5     description='Sample ETL DAG using Bash',
6     schedule_interval=timedelta(days=1),
7 )
```

Here we are creating a variable named `dag` by instantiating the `DAG` class with the following parameters.

`sample-etl-dag` is the ID of the DAG. This is what you see on the web console.

We are passing the dictionary `default_args`, in which all the defaults are defined.

`description` helps us in understanding what this DAG does.

`schedule_interval` tells us how frequently this DAG runs. In this case every day. (`days=1`).

Exercise 5 - Submit a DAG

Submitting a DAG is as simple as copying the DAG python file into `dags` folder in the `AIRFLOW_HOME` directory.

Open a terminal and run the command below to submit the DAG that was created in the previous exercise.

Note: While submitting the dag that was created in the previous exercise, use `sudo` in the terminal before the command used to submit the dag.

```
1 cp my_first_dag.py $AIRFLOW_HOME/dags
```

Next, run the command below one by one to submit shell script in the dags folder and to change the permission for reading shell script.

```
1 cp my_first_dag.sh $AIRFLOW_HOME/dags  
2 cd airflow/dags  
3 chmod 777 my_first_dag.sh
```

Verify that our DAG actually got submitted.

Run the command below to list out all the existing DAGs.

```
1 airflow dags list
```

Verify that `my-first-dag` is a part of the output.

```
1 airflow dags list |grep "my-first-dag"
```

You should see your DAG name in the output.

Run the command below to list out all the tasks in `my-first-dag`.

```
1 airflow tasks list my-first-dag
```

You should see 1 task in the output.

Airflow Monitoring and Logging

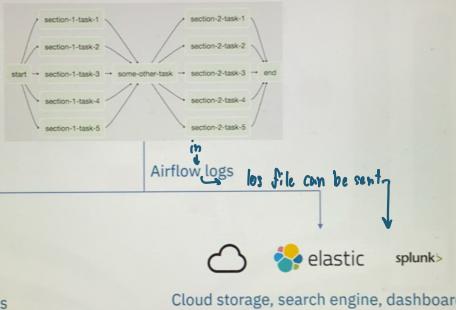
Objectives

After watching this video, you will be able to:

- Use logging capabilities to monitor the task status and diagnose problems with DAG runs
- Explain accessing emitted metrics such as counters, gauges, and timers

Logging

by default, Airflow logs are saved to local file system as log files



Airflow log files

Default log file location:

log file are organized by dag id

logs/dag_id/task_id/execution_date/try_number.log if you want to find the log
↓ of the first execution try
of task 2 in
a dummy DAG

logs/dummy_dag/task1/2021-07-29T00:17:00+00:00/1.log
in just log file

[2021-07-29 17:29:09,778] {subprocess.py:3} INFO - Running command: ['bash', '-c', 'sleep 1'] ← Permissions COMMAND
[2021-07-29 17:29:09,778] {subprocess.py:3} INFO - Output: ← Command result
[2021-07-29 17:29:09,851] {subprocess.py:3} INFO - Command exited with return code 0 ← Command result
[2021-07-29 17:29:09,777] {taskinstance.py:1103} INFO - Marking task as SUCCESS, dag_id=dummy_dag, task_id=task1, execution_date=2021-07-29T00:00:00, start_date=2021-07-29T17:29:00
[2021-07-29 17:29:09,777] {taskinstance.py:1124} INFO - 0 downstream tasks scheduled from follow-on schedule check
[2021-07-29 17:29:09,879] {local_task_job.py:510} INFO - Task exited with return code 0 ← Task result

Reviewing task events via UI

ID	Dttm	Dag Id	Task Id	Event	Execution Date	Owner	Extra
613	2021-07-29, 17:32:23	dummy_dag	task1	success	2021-07-29, 00:51:00	Rameesh Sannareddy	
604	2021-07-29, 17:32:20	dummy_dag	task1	ctc_task_run	2021-07-29, 00:51:00	default	{"host_name": "2550fa70ba", "full_command": ["home/airflow/local/bin/airflow", "celery", "worker"]}
603	2021-07-29, 17:32:20	dummy_dag	task1	running	2021-07-29, 00:51:00	Rameesh Sannareddy	
602	2021-07-29, 17:32:19	dummy_dag	task1	success	2021-07-29, 00:50:00	Rameesh Sannareddy	
598	2021-07-29, 17:32:16	dummy_dag	task1	ctc_task_run	2021-07-29, 00:51:00	default	{"host_name": "2550fa70ba", "full_command": ["home/airflow/local/bin/airflow", "celery", "worker"]}
594	2021-07-29, 17:32:15	dummy_dag	task1	ctc_task_run	2021-07-29, 00:50:00	default	{"host_name": "2550fa70ba", "full_command": ["home/airflow/local/bin/airflow", "celery", "worker"]}

Monitoring metrics

Airflow produces three different types of metrics for checking

- **Counters:** Metrics that always increase
 - Total count of task instances failures
 - Total count of task instances successes
- **Gauges:** Metrics that may fluctuate
 - Number of running tasks
 - DAG bag size, or number of DAGs in production
- **Timers:** Metrics related to time duration
 - Milliseconds to finish a task
 - Milliseconds to reach a success or failed state

Storing and analyzing metrics



