Hands-on Lab: Create a DAG for Apache Airflow with BashOperator



Estimated time needed: 40 minutes

Introduction

In this lab, you will create workflows using BashOperator in Airflow DAGs and simulate an ETL process using bash commands that are scheduled to run once a day.

Objectives

After completing this lab, you will be able to:

- Explore the Airflow Web UI
- Create a DAG with BashOperator
- Submit a DAG and run it through Web UI

Prerequisites

Please ensure that you have completed the reading on the Airflow DAG Operators before proceeding with this lab. It is highly recommended that you are familiar with bash commands to do this lab.

About Skills Network Cloud IDE

Skills Network Cloud IDE (based on Theia and Docker) provides an environment for hands-on labs for course and project-related labs. Theia is an open-source IDE (Integrated Development Environment) that can be run on a desktop or on the cloud. To complete this lab, you will be using the Cloud IDE based on Theia, running in a Docker container.

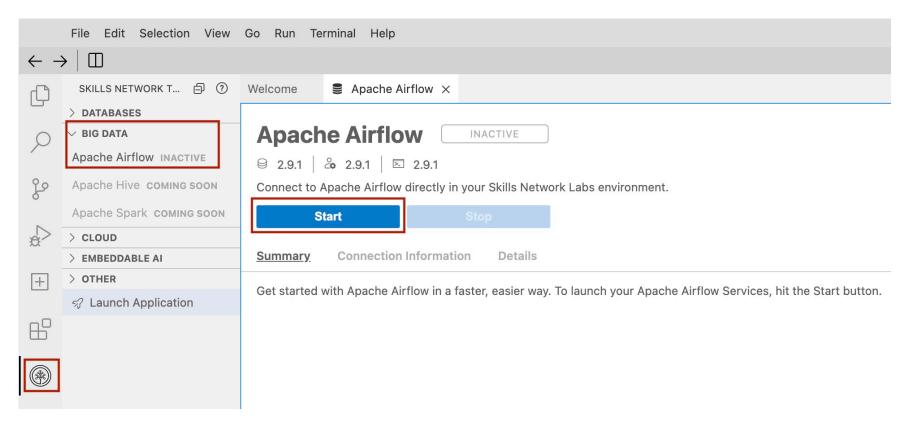
Important notice about this lab environment

Please be aware that sessions for this lab environment are not persistent. A new environment is created for you every time you connect to this lab. Any data you may have saved in an earlier session will get lost. To avoid losing your data, please plan to complete these labs in a single session.

Exercise 1: Start Apache Airflow

- 1. Click on Skills Network Toolbox.
- 2. From the BIG DATA section, click Apache Airflow.
- 3. Click **Start** to start the Apache Airflow.

about:blank 1/10

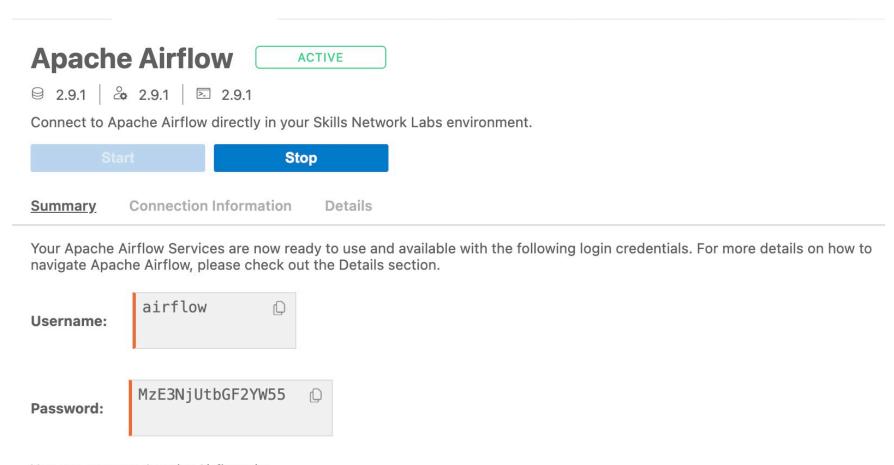


Note: Please be patient, it will take a few minutes for Airflow to start.

Exercise 2: Open the Airflow Web UI

1. When Airflow starts successfully, you should see an output similar to the one below. Once Apache Airflow has started, click on the highlighted icon to open Apache Airflow Web UI in the new window.

about:blank 2/10

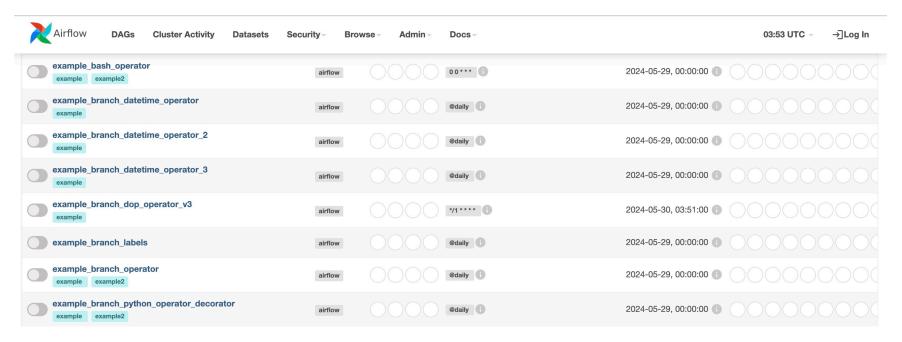


You can manage Apache Airflow via:



You should land on a page that looks like this.

about:blank 3/10



Exercise 3: Create a DAG

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

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

- 1. .
- 2. 2 3. 3
- 1 1
- 5. 5 6. 6
- 7. 7
- 8.8
- 9. 9 10. 10
- 11. 11
- 12. 12
- 13. 13
- 14. 14
- 15. 15 16. 16
- 17. 17
- 18. 18
- 19. 19 20. 20
- 21. 21 22. 22
- 23. 23
- 24. 24
- 25. 25
- 26. 26 27. 27
- 28. 28
- 29. 29
- 30. 30

about:blank

```
11/9/24, 3:29 PM
    31. 31
   32. 32
   33. 33
    34. 34
   35. 35
    36. 36
    37. 37
   38. 38
    39. 39
    40. 40
    41. 41
    42. 42
    43. 43
   44. 44
   45. 45
    46.46
    47. 47
    48. 48
    49. 49
    50. 50
    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 operator initialization
   14. default_args = {
            'owner': 'your_name_here',
'start_date': days_ago(0),
   15.
   16.
            'email': ['your_email_here'],
   17.
            'retries': 1,
   18.
            'retry_delay': timedelta(minutes=5),
   19.
   20. }
    21.
    22. # defining the DAG
   23.
   24. # define the DAG
    25. dag = DAG(
            'my-first-dag',
   26.
            default_args=default_args,
   27.
            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.
            bash_command='cut_-d":" -f1,3,6 /etc/passwd > /home/project/airflow/dags/extracted-data.txt',
    38.
   39.
            dag=dag,
    40.)
   41.
    42. # define the second task
    43. transform and load = BashOperator(
            task_id='transform',
    44.
            bash_command='tr ":" "," < /home/project/airflow/dags/extracted-data.txt > /home/project/airflow/dags/transformed-data.csv',
    45.
    46.
            dag=dag,
    47.)
    48.
    49. # task pipeline
    50. extract >> transform_and_load
```

Copied!

- 1. Create a new file by choosing File->New File and naming it my_first_dag.py.
- 2. Then, copy the code above and paste it into my_first_dag.py.

Exercise 4: Submit a DAG

Submitting a DAG is as simple as copying the DAG Python file into the 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 airflow 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.

1. Open a terminal and run the command below to set the AIRFLOW HOME.

```
1. 1
2. 2
1. export AIRFLOW_HOME=/home/project/airflow
2. echo $AIRFLOW_HOME
Copied!
```

```
theia@theiadocker-lavanyas:/home/project ×

theia@theiadocker-lavanyas:/home/project$ echo $AIRFLOW_HOME
/home/project/airflow
```

2. Run the command below to submit the DAG that was created in the previous exercise.

```
1. 1
2. 2
1. export AIRFLOW_HOME=/home/project/airflow
2. cp my_first_dag.py $AIRFLOW_HOME/dags
Copied!
```

3. Verify that your DAG actually got submitted.

about:blank 6/10

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

```
1. 1
1. airflow dags list
Copied!
```

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

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

You should see your DAG name in the output.

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

```
1. 1
1. airflow tasks list my-first-dag
Copied!
```

You should see 2 tasks in the output.

Practice exercise

Write a DAG named ETL_Server_Access_Log_Processing.py.

- 1. Create the imports block.
- 2. Create the DAG Arguments block. You can use the default settings
- 3. Create the DAG definition block. The DAG should run daily.
- 4. Create the download task. The download task must download the server access log file, which is available at the URL:

https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DB0250EN-SkillsNetwork/labs/Apache%20Airflow/Build%20a%20DAG%20using%20Airflow/web-server-access-log.txt

5. Create the extract task.

The server access log file contains these fields.

```
a. timestamp - TIMESTAMP
b. latitude - float
c. longitude - float
d. visitorid - char(37)
e. accessed_from_mobile - boolean
f. browser_code - int
```

The extract task must extract the fields timestamp and visitorid.

- 6. Create the transform task. The transform task must capitalize the visitorid.
- 7. Create the load task. The load task must compress the extracted and transformed data.
- 8. Create the task pipeline block. The pipeline block should schedule the task in the order listed below:
 - 1. download
 - extract
 - 3. transform
 - 4. load
- 9. Submit the DAG.
- 10. Verify if the DAG is submitted.
- ► Click here for **hint**.
- ▼ Click here for the **solution**.

Add to the file the following parts of code to ETL_Server_Access_Log_Processing.py to complete the tasks given in the problem.

- 1. 1 2. 2 3. 3
- 4. 4 5. 5
- 6.6
- 7. 7
- 8.8 9.9
- 10. 10
- 11. 11 12. 12
- 13. 13
- 14. 14
- 15. 15
- 16. 16
- 17. 17 18. 18
- 19. 19
- 20. 20 21. 21
- 22. 22 23. 23 24. 24
- 25. 25
- 26. 26 27. 27
- 28. 28
- 29. 29 30. 30
- 31. 31 32. 32 33. 33

- 33. 33 34. 34 35. 35 36. 36 37. 37
- 38. 38 39. 39
- 40. 40
- 41. 41 42. 42
- 43. 43
- 44. 44 45. 45
- 46. 46
- 47. 47
- 48. 48
- 49. 49
- 50. 50 51. 51 52. 52
- 53. 53 54. 54 55. 55
- 56. 56
- 57. 57 58. 58
- 59. 59
- 60. 60 61. 61
- 62. 62
- 63. 63 64. 64
- 65. 65
- 66. 66 67. 67
- 68. 68
- 69. 69
- 1. # import the libraries

about:blank

```
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 operator initialization
14. default args = {
         'owner': 'your_name',
15.
         'start_date': days_ago(0),
16.
17.
         'email': ['your email'],
18.
         'retries': 1,
19.
         'retry_delay': timedelta(minutes=5),
20. }
21.
22. # defining the DAG
23.
24. # define the DAG
25. dag = DAG(
26.
         'ETL_Server_Access_Log_Processing',
27.
         default args=default args,
28.
        description='My first DAG',
         schedule interval=timedelta(days=1),
30.)
31.
32. # define the tasks
33.
34. # define the task 'download'
35.
36. download = BashOperator(
37.
        task id='download',
38.
        bash_command='curl "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DB0250EN-SkillsNetwork/labs/Apache%20Airflow/Build%20a%20DAG%20using%20Airflow/web-server-access-log.txt
39.
40.)
41.
42. # define the task 'extract'
43.
44. extract = BashOperator(
45.
         task id='extract',
46.
        bash_command='cut -f1,4 -d"#" web-server-access-log.txt > /home/project/airflow/dags/extracted.txt',
47.
        dag=dag,
48.)
49.
50.
51. # define the task 'transform'
52.
53. transform = BashOperator(
54.
        task id='transform',
55.
        bash_command='tr "[a-z]" "[A-Z]" < /home/project/airflow/dags/extracted.txt > /home/project/airflow/dags/capitalized.txt',
56.
        dag=dag,
57.)
58.
59. # define the task 'load'
61. load = BashOperator(
62.
        task id='load',
63.
        bash_command='zip log.zip capitalized.txt' ,
64.
65.)
66.
67. # task pipeline
69. download >> extract >> transform >> load
Copied!
```

Submit the DAG by running the following command.

about:blank 9/10

1. 1

cp ETL_Server_Access_Log_Processing.py \$AIRFLOW_HOME/dags



Verify if the DAG is submitted on the Web UI or the CLI using the below command.

1. airflow dags list



Authors

<u>Lavanya T S</u> Ramesh Sannareddy

© IBM Corporation. All rights reserved.

about:blank 10/10