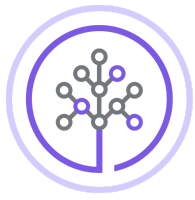


Hands-on Lab: Create and execute a Shell script from Airflow



Skills
Network

Estimated time needed: **40** minutes

Objectives

After completing this lab you will be able to:

- Create a basic shell script
- Explore the anatomy of a DAG.
- Create a DAG.
- Call and execute the shell script
- Submit a DAG.

About Skills Network Cloud IDE

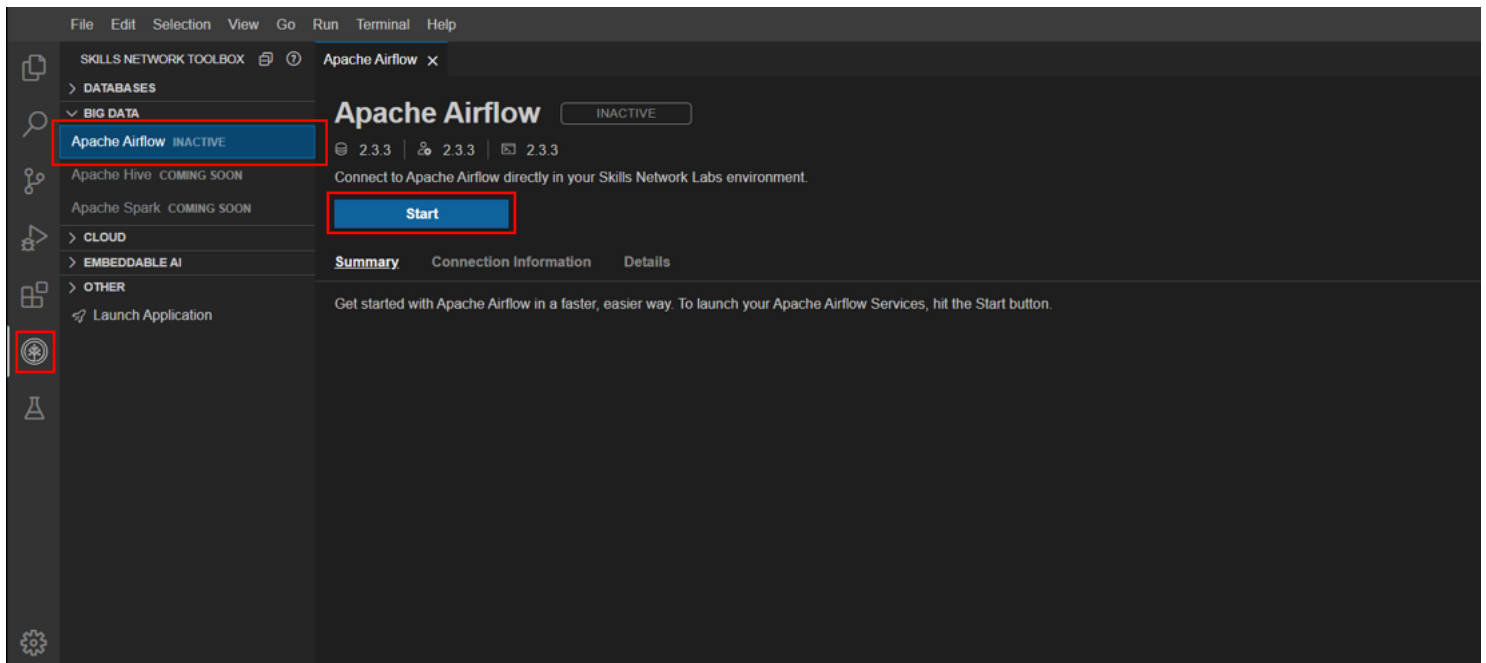
Skills Network Cloud IDE (based on Theia and Docker) provides an environment for hands on labs for course and project related labs. Theia is an open source IDE (Integrated Development Environment), that can be run on desktop or on the cloud. to complete this lab, we will be using the Cloud IDE based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

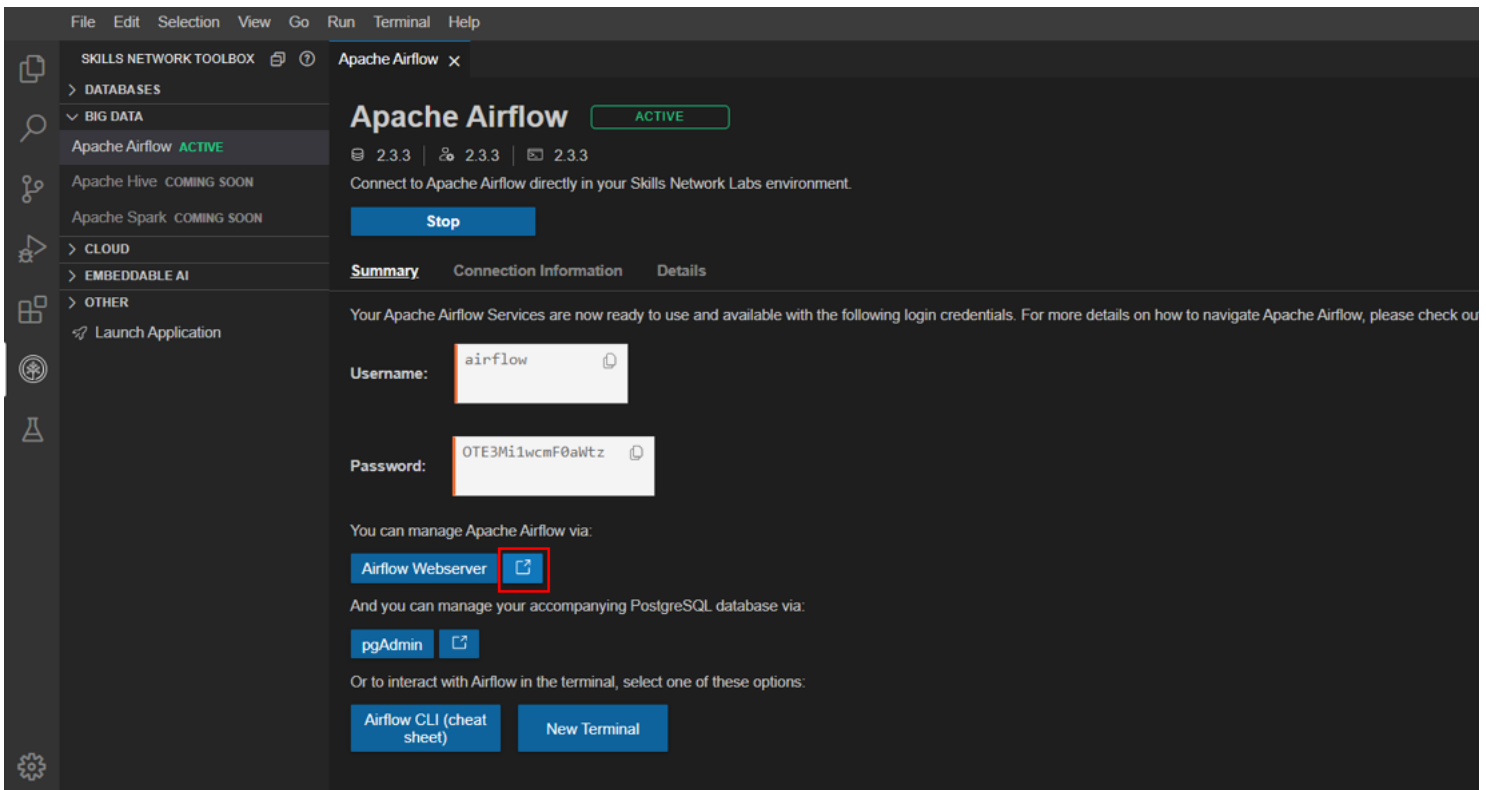
Click on **Skills Network Toolbox**. In **BIG DATA** section, click **Apache Airflow**. To start the Apache Airflow click **Start**.



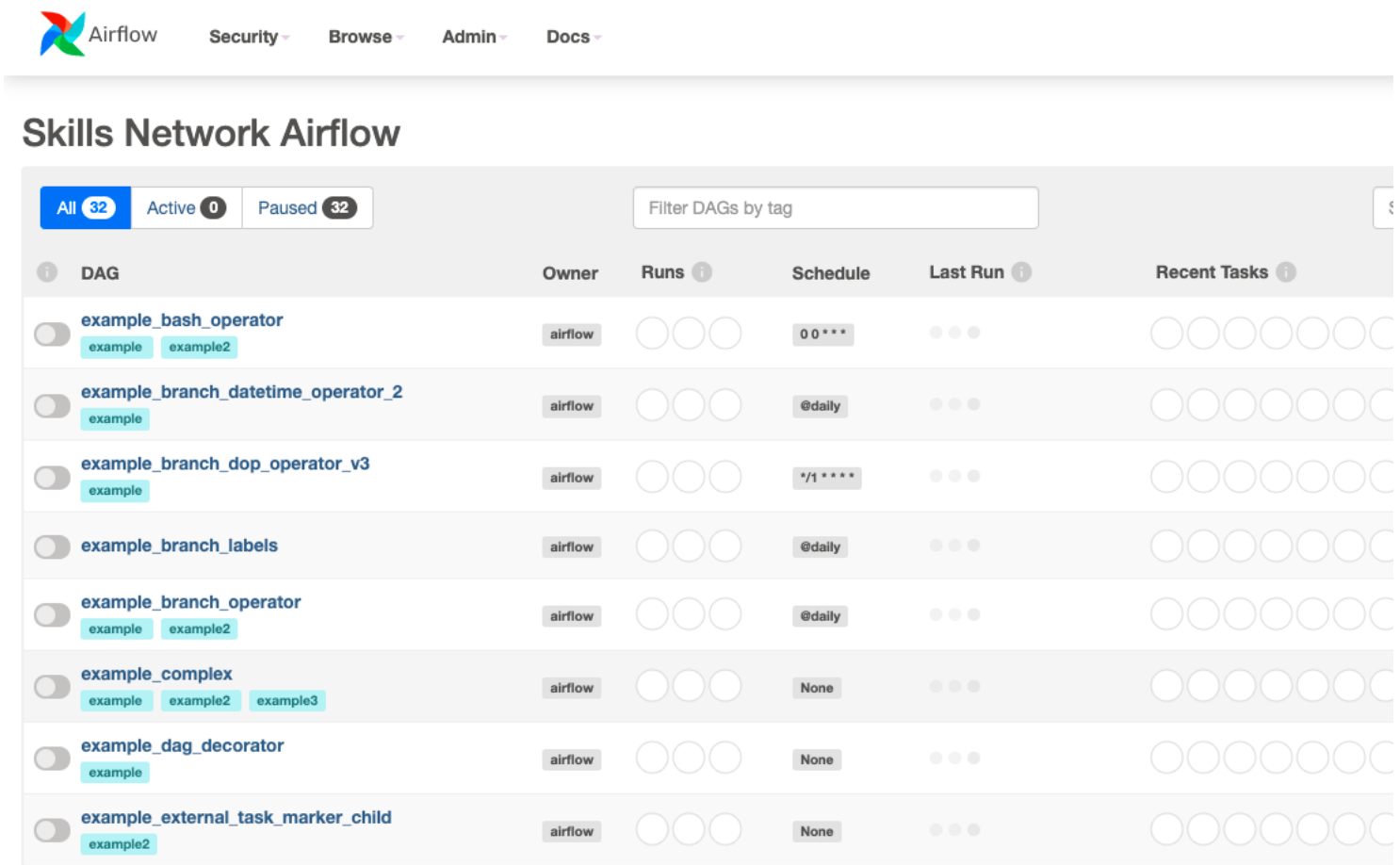
Please be patient, it will take a few minutes for airflow to get started.

Exercise 2 - Open the Airflow Web UI

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.



You should land at a page that looks like this.



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. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7

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

Copied!

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

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. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9

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
```

Copied!

A typical DAG Arguments block looks like this.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12

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@somemail.com'],
8.     'email_on_failure': True,
9.     'email_on_retry': True,
10.    'retries': 1,
11.    'retry_delay': timedelta(minutes=5),
12. }
```

Copied!

DAG arguments are like settings for the DAG.

The above settings mention

- the owner name,
- when this DAG should run from: `days_age(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.

A typical DAG definition block looks like this.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7

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. )
```

Copied!

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`).

A typical task definitions block looks like this:

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8

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. )
```

Copied!

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. 1
2. 2

1. # task pipeline
2. extract_transform_and_load
```

Copied!

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

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. 1
2. 2
3. 3
4. 4
5. 5

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
```

Copied!

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 `et1` that calls the shell script `my_first_dag.sh`

```
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
34. 34
35. 35
36. 36
37. 37
38. 38
39. 39
40. 40
41. 41
42. 42
43. 43
44. 44

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@somemail.com'],
18.     'email_on_failure': False,
19.     'email_on_retry': False,
20.     'retries': 1,
21.     'retry_delay': timedelta(minutes=5),
22. }
23.
```

```

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

```

Copied!

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

```

Copied!

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. 1
2. 2
3. 3

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

```

Copied!

Verify that our DAG actually got submitted.

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

```

1. 1
1. airflow dags list

```

Copied!

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.

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 1 task in the output.

Practice exercises

1. Problem:

Download the dataset from the source to the destination mentioned below using `wget` command in terminal.

Note: While downloading the file in the terminal use the **sudo** command before the command used to download the file.

Source : <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>

Destination : /home/project/airflow/dags.

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

Write a shell script named `ETL_Server_Access_Log_Processing.sh`.

Task 1: Add a command in the shell script to extract the fields `timestamp` and `visitorid` from the `web-server-access-log.txt` and write to a file `extracted.txt`

Task 2: Update the shell script to add a command to capitalize the `visitorid` and write to a new file `capitalized.txt`.

Task 3: Update the shell script to add a command to compress the transformed file `capitalized.txt` into a new file `log.tar.gz`.

Write a DAG named `ETL_Server_Access_Log_Processing`.

Task 2: Create the imports block.

Task 3: Create the DAG Arguments block. You can use the default settings

Task 4: Create the DAG definition block. The DAG should run daily.

Task 5: Create the task `extract_transform_and_load` to call the shell script.

Task 6: Create the task pipeline block.

Task 7: Submit the DAG.

Task 8: Submit the shell script to `dags` folder.

Task 9: Change the permission to read shell script.

Task 10. Verify if the DAG is submitted

▼ Click here for Hint

Follow the example python code given in the lab and make necessary changes to create the new DAG.

▼ Click here for Solution

Download the dataset from the source to the destination mentioned below using `wget` command in terminal.

Note: While downloading the file in the terminal use the **sudo** command before the command used to download the file.

Source : <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>

Destination : /home/project/airflow/dags

Select **File -> New File** from the menu and name it as `ETL_Server_Access_Log_Processing.sh`.

Copy the code below in the shell script.

Task 1: Create a shell script having the following commands.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10

1. #!/bin/bash
2. echo "extract_transform_load"
3. # cut command to extract the fields timestamp and visitorid writes to a new file extracted.txt
4. cut -f1,4 -d"#" /home/project/airflow/dags/web-server-access-log.txt > /home/project/airflow/dags/extracted.txt
5.
6. # tr command to transform by capitalizing the visitorid.
7. tr "[a-z]" "[A-Z]" < /home/project/airflow/dags/extracted.txt > /home/project/airflow/dags/capitalized.txt
8.
```

```
9. # tar command to compress the extracted and transformed data.
10. tar -czvf /home/project/airflow/dags/log.tar.gz /home/project/airflow/dags/capitalized.txt
```

Copied!

Next select **File -> New File** from the menu and name it as `ETL_Server_Access_Log_Processing.py`.

Task 2: Create the imports block.

Copy the code below in the python file

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9

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
```

Copied!

Task 3: Create the DAG Arguments block. You can use the default settings.

Copy the code below in the python file

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10
11. 11
12. 12

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@somemail.com'],
8.     'email_on_failure': False,
9.     'email_on_retry': False,
10.    'retries': 1,
11.    'retry_delay': timedelta(minutes=5),
12. }
```

Copied!

Task 4: Create the DAG definition block. The DAG should run daily.

Copy the code below in the python file

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9

1. # defining the DAG
2.
3. # define the DAG
4. dag = DAG(
5.     'ETL_Server_Access_Log_Processing',
6.     default_args=default_args,
7.     description='My first DAG',
8.     schedule_interval=timedelta(days=1),
9. )
```

Copied!

Task 5: Create the task `extract_transform_and_load` to call the shell script.

Copy the code below in the python file


```

1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
8. 8
9. 9
10. 10

1. # define the tasks
2.
3. #define the task named extract_transform_and_load to call the shell script
4. #calling the shell script
5. extract_transform_and_load = BashOperator(
6.     task_id="extract_transform_and_load",
7.     bash_command="/home/project/airflow/dags/ETL_Server_Access_Log_Processing.sh ",
8.     dag=dag,
9. )
10.

```

Copied!

Task 6: Create the task pipeline block.

Copy the code below in the python file

```

1. 1
2. 2
3. 3

1. # task pipeline
2.
3. extract_transform_and_load

```

Copied!

Save the python file

Task 7: Submit the DAG.

```

1. 1

1. cp ETL_Server_Access_Log_Processing.py $AIRFLOW_HOME/dags

```

Copied!

Task 8: Submit the shell script to dags folder.

```

1. 1

1. cp ETL_Server_Access_Log_Processing.sh $AIRFLOW_HOME/dags

```

Copied!

Task 9: Change the permission to read shell script.

```

1. 1
2. 2

1. cd airflow/dags
2. chmod 777 ETL_Server_Access_Log_Processing.sh

```

Copied!

Task 10: Verify if the DAG is submitted.

```

1. 1

1. airflow dags list

```

Copied!

Verify that the DAG's Python script `ETL_Server_Access_Log_Processing.py` is listed.

Authors

Ramesh Sannareddy

Pratiksha Verma

Other Contributors

Rav Ahuja

Lakshmi Holla

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2023-03-02	0.6	Pratiksha Verma	Modified instructions to create DAG using shell script
2022-11-10	0.5	Appalabhaktula Hema	Updated instruction
2022-08-22	0.4	Lakshmi Holla	updated bash command
2022-07-29	0.3	Lakshmi Holla	changed dag name
2022-06-28	0.2	Lakshmi Holla	updated DAG path
2021-07-05	0.1	Ramesh Sannareddy	Created initial version of the lab

Copyright (c) 2021 IBM Corporation. All rights reserved.