

# Final Demonstration

This document will guide you through the demonstration of the real-world use case of Airflow in Session 3 of the Airflow module.

## Prerequisites:

- mysql.dump (MySQL queries for table creation and record )
- etl\_dag.py(the code for the DAG explained in the video)
- filter\_trip.py (Spark application)
- generate\_trip\_throughput.py (Spark application)
- filter\_booking.py (Spark application)
- generate\_car\_with\_most\_trips.py (Spark application)
- Make sure that the current Java version is Java 11(11.x.x). **Note:** This is only needed to run the Sqoop operator and you will need to switch back to Java 8 after you are done with the Sqoop operator. The steps to switch back to Java 8 can be found at the end of this document.

You can check this by running the following command:

```
java -version
```

```
(airflow) [hadoop@ip-172-31-34-198 ~]$ java -version
openjdk version "11.0.13" 2021-10-19 LTS
OpenJDK Runtime Environment 18.9 (build 11.0.13+8-LTS)
OpenJDK 64-Bit Server VM 18.9 (build 11.0.13+8-LTS, mixed mode, sharing)
```

If you still have Java 8, then you need to switch to Java 11 by running the following command.

```
sudo alternatives --config java <<< 3
```

```
(airflow) [hadoop@ip-172-31-33-82 ~]$ sudo alternatives --config java <<< 3
There are 3 programs which provide 'java'.

  Selection    Command
-----
  1            /usr/lib/jvm/java-1.8.0-amazon-corretto.x86_64/jre/bin/java
*  2            /usr/lib/jvm/java-17-amazon-corretto.x86_64/bin/java
+  3            java-11-openjdk.x86_64 (/usr/lib/jvm/java-11-openjdk-11.0.13.0.8-1.amzn2.0.3.x86_64/bin/java)

Enter to keep the current selection[+], or type selection number: (airflow) [hadoop@ip-172-31-33-82 ~]$ |
```

### What are we doing?

We will be creating a DAG for the ride-hailing problem statement explained in Session 3

### Please follow the instructions below:

1. Login to your EMR instance.
2. Activate the Python virtual environment using the following command:

**source /home/hadoop/airflow/bin/activate**

```
[hadoop@ip-172-31-33-82 ~]$ source /home/hadoop/airflow/bin/activate  
(airflow) [hadoop@ip-172-31-33-82 ~]$ |
```

3. Now you need to load the data into your local MySQL, firstly you need to place the mysql.dump file in some location in your EMR machine. We will store in the /home/hadoop/ directory

(You can use WinSCP or create a new file called mysql.dump in the /home/hadoop/ directory and paste the contents in that file)

4. Now run the following command to execute the SQL commands in /tmp/transactions.dump and create our tables :

**mysql -u root -p123 < /home/hadoop/mysql.dump**

5. Next, you need to set up the different connections using the Airflow UI which is hosted in the URL : **your\_public\_dns:8082**

**Note:** You can find your\_public\_ip in your AWS EMR dashboard (IPv4 Public DNS))

Edit the following connections:

#### **Sqoop:**

Conn Id: sqoop\_default

Conn Type: Sqoop (Select from the drop-down )

Connection URL: **jdbc:mysql://<public DNS>**

Schema: **events**

Login: **root**

Password: **123**

Connection Id *	sqoop_default
Connection Type *	Sqoop
Connection Type missing? Make sure you've installed the corresponding client.	
Description	
Host	jdbc:mysql://ec2-34-201-68-160.compute-1.amazonaws.com
Schema	events
Login	root

#### Hive:

Conn Id: hive\_cli\_default

Conn Type: Hive Client Wrapper (Select from the drop-down )

Host: **<private IP of master node>**

Login: **hadoop**

Port: 10000

Extra: {"use\_beeline": true, "auth": ""}

Connection Id *	hive_cli_default
Connection Type *	<div>Hive Client Wrapper</div> <div>Connection Type missing? Make sure you've installed the corresponding Airflow Provider Package.</div>
Description	
Host	172.31.58.178
Schema	default
Login	hadoop
Password	
Port	10000
Extra	<pre>{"use_beeline": true, "auth": ""}</pre>

### Spark:

Conn Id: spark\_default

Conn Type: Spark (Select from the drop-down )

Host: yarn

Extra: **`{"master": "yarn", "conf": "/etc/spark/conf/spark-defaults.conf"}`**

Connection Id *	spark_default
Connection Type *	Spark Connection Type missing? Make sure you've installed the correspond
Description	
Host	yarn
Port	
Extra	{"master": "yarn", "conf": "/etc/spark/conf/spark-defaults.conf"}

Save Test

- Create a directory called **uber** in the `airflow_codes` directory using the following command:

```
mkdir -p /home/hadoop/airflow_codes/uber
```

- Now place the following file inside the **uber directory** you just created:

- filter\_trip.py
- generate\_trip\_throughput.py
- filter\_booking.py
- generate\_car\_with\_most\_trips.py

(You can use WinSCP or create a new file and paste the code in that file)

```
(airflow) [hadoop@ip-172-31-58-178 ~]$ ls airflow_codes/uber/
filter_booking.py  generate_car_with_most_trips.py
filter_trip.py     generate_trip_throughput.py
(airflow) [hadoop@ip-172-31-58-178 ~]$ |
```

8. Now you need to place the **etl\_dag.py** file in the **/home/hadoop/airflow/dags** directory.  
(You can use WinSCP or create a new file and paste the code in that file)
9. To ensure that the file there are no issues/errors with the file is it considered good practice to compile the program using the following command:

**python etl\_dag.py**

10. You can also use the following command to list the dags in your instance:

**airflow dags list**

11. Once you have made sure that your dag file has no issues you can go back to the Airflow UI
12. In case you are re-running this DAG, you will have to delete the target\_dir of the sqoop task as we did in the SqoopOperator demonstration to avoid any errors.

You can enter the following command to do so:

**sudo hdfs dfs -rm -r -skipTrash /data/raw**

Also, clear the task/DAG before re-running it.

13. Switch ON the DAG(etl\_dag)

<div> <div>All 7</div> <div>Active 1</div> <div>Paused 6</div> <div>Filter</div> </div>			
DAG	Owner	Runs	Schedule
<input checked="" type="checkbox"/> etl_dag	admin	<div> <div></div> <div>1</div> <div></div> <div>1</div> </div>	@once
<input type="checkbox"/> etl_dag_refined	admin	<div> <div></div> <div></div> <div></div> <div></div> </div>	@once

(Note: The DAG might take a while to show up on the UI. Keep refreshing and wait patiently)

14. Click on the etl\_dag and go to the graph view

You will see the task is running

Tree

Graph

Calendar

Task Duration

Task Trie



2022-02-21T20:08:47Z

Runs

25



Run

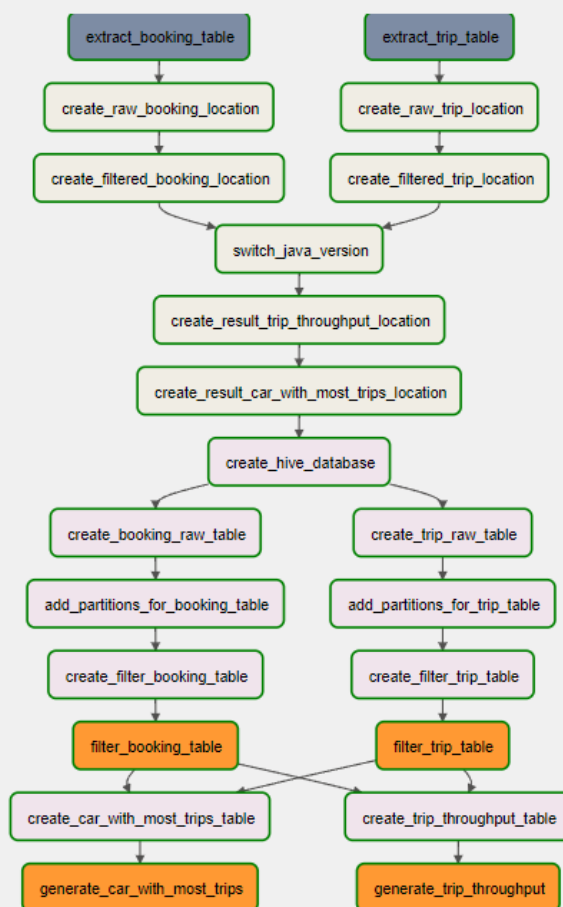
schedule

BashOperator

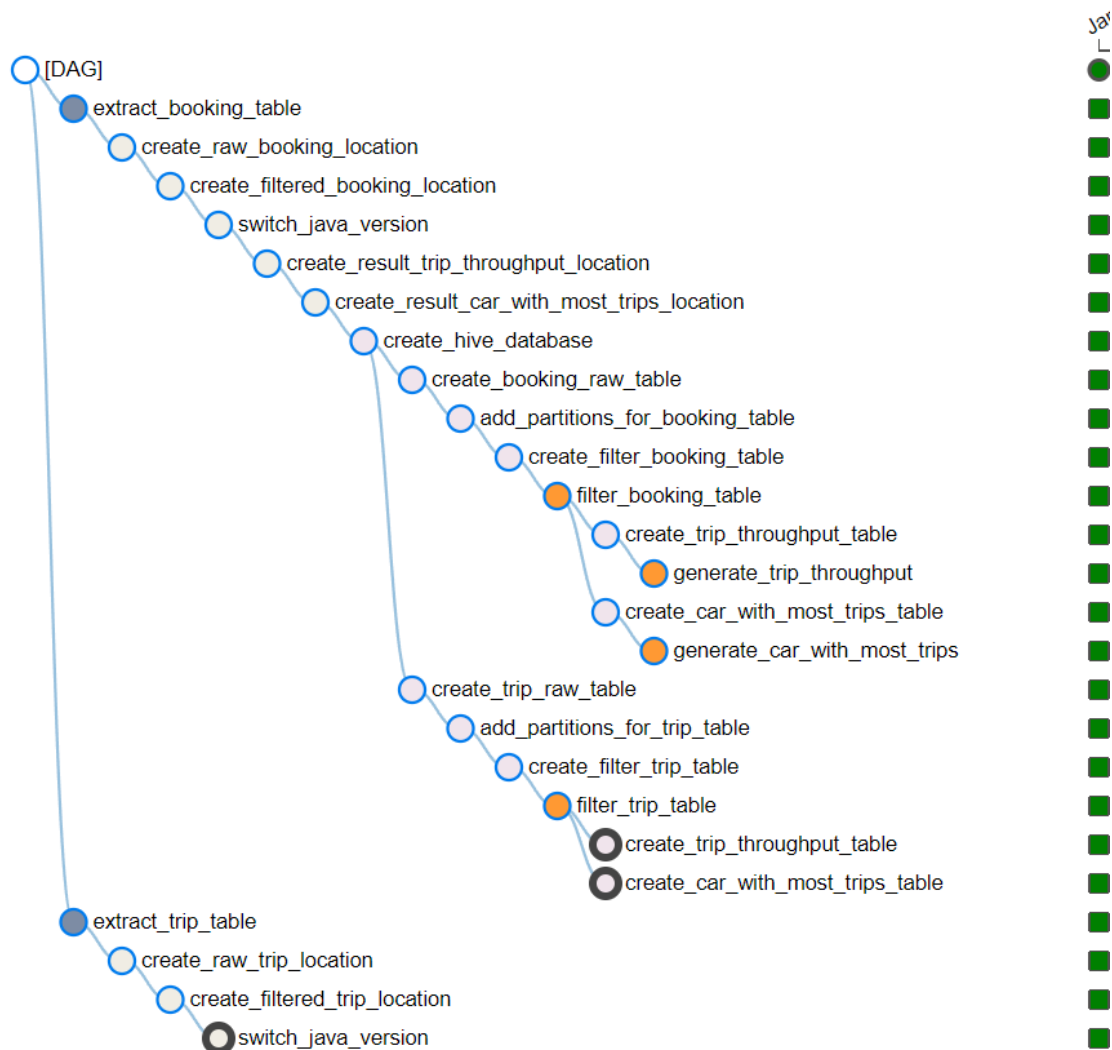
HiveOperator

SparkSubmitOperator

SqoopOperator



This may take a while. So wait patiently.



Click on refresh and eventually, all tasks will have successfully completed

15. Once the DAG has completed execution, the output will be generated in the tables **trip\_troughput** and **car\_with\_most\_trips** inside the **events** database

16. You can view the results in the CLI by following the steps below:

For trip\_throughput enter command:

**hive -e "select \* from events.trip\_throughput;"**



```
mumbai 0.75      2022-01-16
bangalore      0.9      2022-01-16
chennai 0.8      2022-01-16
Time taken: 6.207 seconds, Fetched: 3 row(s)
```

For car\_with\_most\_trips enter command:

```
hive -e "select * from events.car_with_most_trips ;"
```

```
bangalore      sedan      5      2022-01-16
mumbai economy 3      2022-01-16
mumbai sedan   3      2022-01-16
chennai sedan   5      2022-01-16
Time taken: 5.442 seconds, Fetched: 4 row(s)
```

17. You can switch off your DAG if you don't want it to run anymore.