

EECS E6893 Big Data Analytics

Homework Assignment 4

Submitted by:
Shivam Ojha
UNI: so2639

Task 1: (Helloworld)

1.1 Please find below the screenshots for the installation of Airflow on a virtual machine on Google Cloud Platform.

1. Terminal after successfully starting the webserver

```
(airflow) so2639@hw4-airflow:~$ airflow webserver --port 8080

[2021-11-27 21:26:29,988] (dagbag.py:500) INFO - Filling up the DagBag from /dev/null
[2021-11-27 21:26:29,415] (manager.py:512) WARNING - Refused to delete permission view, assoc with role exists DAG Runs.can_create Admin
Running the Gunicorn Server with:
Workers: 4 sync
Host: 0.0.0.0:8080
Timeout: 120
Logfiles: -
Access Logformat:

[2021-11-27 21:26:32+0000] [11206] [INFO] Starting gunicorn 20.1.0
[2021-11-27 21:26:32+0000] [11206] [INFO] Listening at: http://0.0.0.0:8080 (11206)
[2021-11-27 21:26:32+0000] [11206] [INFO] Using worker: sync
[2021-11-27 21:26:32+0000] [11208] [INFO] Booting worker with pid: 11208
[2021-11-27 21:26:32+0000] [11209] [INFO] Booting worker with pid: 11209
[2021-11-27 21:26:32+0000] [11211] [INFO] Booting worker with pid: 11211
[2021-11-27 21:26:32+0000] [11210] [INFO] Booting worker with pid: 11210
[2021-11-27 21:26:35,795] (manager.py:512) WARNING - Refused to delete permission view, assoc with role exists DAG Runs.can_create Admin
[2021-11-27 21:26:35,964] (manager.py:512) WARNING - Refused to delete permission view, assoc with role exists DAG Runs.can_create Admin
[2021-11-27 21:26:36,076] (manager.py:512) WARNING - Refused to delete permission view, assoc with role exists DAG Runs.can_create Admin
[2021-11-27 21:26:36,081] (manager.py:512) WARNING - Refused to delete permission view, assoc with role exists DAG Runs.can_create Admin
160.39.11.228 - - [27/Nov/2021:21:26:43 +0000] "GET / HTTP/1.1" 302 217 "-" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/15.1 Safari/605.1.15"
160.39.11.228 - - [27/Nov/2021:21:26:44 +0000] "GET /static/appbuilder/css/font-awesome.min.css HTTP/1.1" 200 0 "http://35.185.10.12:8080/home" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/15.1 Safari/605.1.15"
160.39.11.228 - - [27/Nov/2021:21:26:44 +0000] "GET /static/appbuilder/select2/select2.css HTTP/1.1" 200 0 "http://35.185.10.12:8080/home" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/15.1 Safari/605.1.15"
160.39.11.228 - - [27/Nov/2021:21:26:44 +0000] "GET /static/appbuilder/css/flags/flags16.css HTTP/1.1" 200 0 "http://35.185.10.12:8080/home" "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/15.1 Safari/605.1.15"
```

Terminal after successfully starting the scheduler

```
(airflow) so2639@hw4-airflow:~$ airflow scheduler

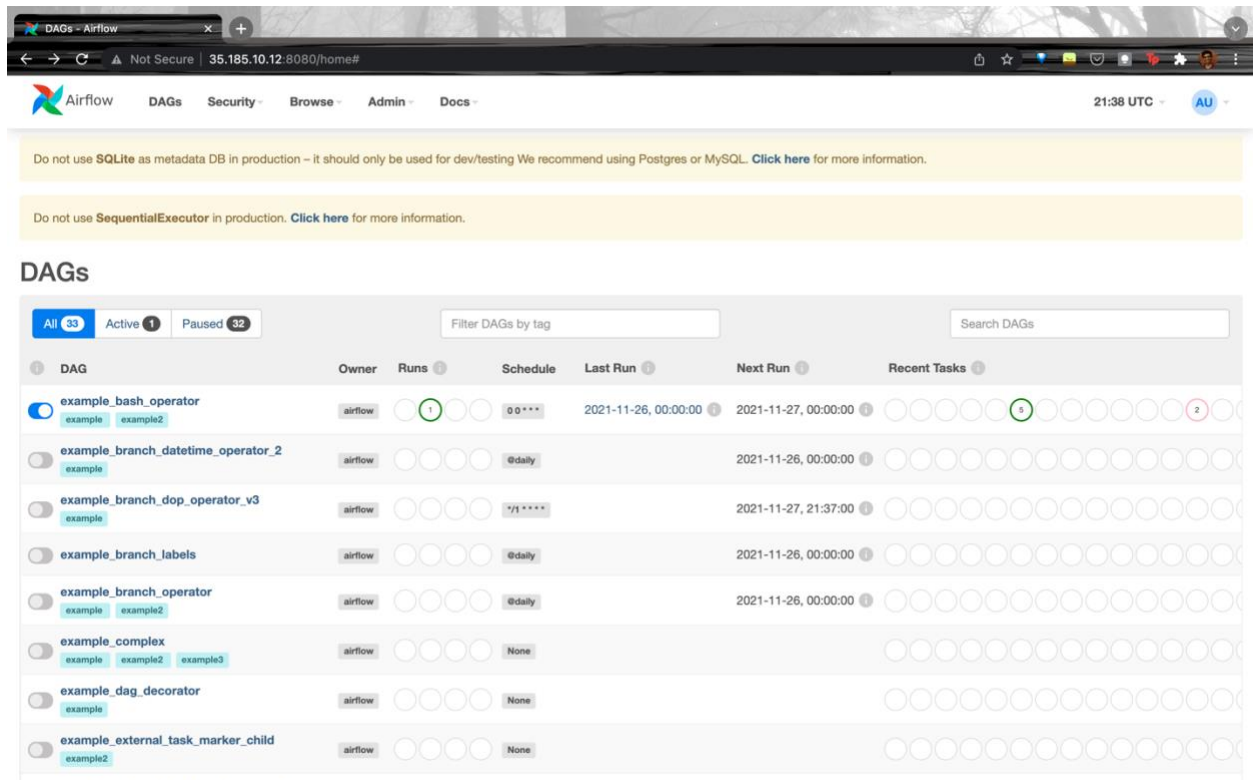
[2021-11-27 21:36:30,696] (scheduler_job.py:596) INFO - Starting the scheduler
[2021-11-27 21:36:30,699] (scheduler_job.py:601) INFO - Processing each file at most -1 times
[2021-11-27 21:36:30+0000] [11464] [INFO] Starting gunicorn 20.1.0
[2021-11-27 21:36:30,785] (manager.py:163) INFO - Launched DagFileProcessorManager with pid: 11465
[2021-11-27 21:36:30+0000] [11464] [INFO] Listening at: http://0.0.0.0:8793 (11464)
[2021-11-27 21:36:30+0000] [11464] [INFO] Using worker: sync
[2021-11-27 21:36:30+0000] (scheduler_job.py:1115) INFO - Resetting orphaned tasks for active dag runs
[2021-11-27 21:36:30+0000] [11466] [INFO] Booting worker with pid: 11466
[2021-11-27 21:36:30,717] (settings.py:52) INFO - Configured default timezone Timezone('UTC')
[2021-11-27 21:36:30,738] (manager.py:431) WARNING - Because we cannot use more than 1 thread (parsing_processes = 2 ) when using sqlite. So we set parallelism to 1.
[2021-11-27 21:36:30+0000] [11467] [INFO] Booting worker with pid: 11467
[2021-11-27 21:41:30,878] (scheduler_job.py:1115) INFO - Resetting orphaned tasks for active dag runs
[2021-11-27 21:46:31,020] (scheduler_job.py:1115) INFO - Resetting orphaned tasks for active dag runs
[2021-11-27 21:51:31,160] (scheduler_job.py:1115) INFO - Resetting orphaned tasks for active dag runs
[2021-11-27 21:56:31,315] (scheduler_job.py:1115) INFO - Resetting orphaned tasks for active dag runs
[2021-11-27 21:59:14,570] (dag.py:2921) INFO - Setting next dagrun for helloworld to 2021-11-27T21:59:09.346568+00:00
[2021-11-27 21:59:14,628] (scheduler_job.py:288) INFO - 1 tasks up for execution:
<TaskInstance: helloworld.tl scheduled_2021-11-26T21:59:09.346568+00:00 [scheduled]>
[2021-11-27 21:59:14,628] (scheduler_job.py:317) INFO - Figuring out tasks to run in Pool(name=default_pool) with 128 open slots and 1 task instances ready to be queued
[2021-11-27 21:59:14,630] (scheduler_job.py:345) INFO - DAG helloworld has 0/16 running and queued tasks
[2021-11-27 21:59:14,630] (scheduler_job.py:410) INFO - Setting the following tasks to queued state:
<TaskInstance: helloworld.tl scheduled_2021-11-26T21:59:09.346568+00:00 [scheduled]>
[2021-11-27 21:59:14,631] (scheduler_job.py:450) INFO - Sending TaskInstanceKey(dag_id='helloworld', task_id='tl', run_id='scheduled_2021-11-26T21:59:09.346568+00:00', try_number=1) to executor with priority 7 and queue default
[2021-11-27 21:59:14,631] (scheduler_job.py:82) INFO - Adding to queue: ['airflow', 'tasks', 'run', 'helloworld', 'tl', 'scheduled_2021-11-26T21:59:09.346568+00:00', '--local', '--subdir', 'DAGS_FOLDER/helloworld.py']
[2021-11-27 21:59:14,631] (sequential_executor.py:59) INFO - Executing command: ['airflow', 'tasks', 'run', 'helloworld', 'tl', 'scheduled_2021-11-26T21:59:09.346568+00:00', '--local', '--subdir', 'DAGS_FOLDER/helloworld.py']
[2021-11-27 21:59:18,140] (dagbag.py:500) INFO - Filling up the DagBag from /home/so2639/airflow/dags/helloworld.py
Running <TaskInstance: helloworld.tl scheduled_2021-11-26T21:59:09.346568+00:00 [queued]> on host hw4-airflow.c.big-data-6893-326522.internal
[2021-11-27 21:59:18,783] (scheduler_job.py:504) INFO - Executor reports execution of helloworld.tl run_id=scheduled_2021-11-26T21:59:09.346568+00:00 exited with status success for try number 1
[2021-11-27 21:59:18,792] (scheduler_job.py:547) INFO - TaskInstance Finished: dag_id=helloworld, task_id=tl, run_id=scheduled_2021-11-26T21:59:09.346568+00:00, run_start_date=2021-11-27 21:59:16.285079+00:00, run_end_date=2021-11-27 21:59:18.389010+00:00, run_duration=2.103931, state=success, executor_state=success, try_number=1, max_tries=1, job_id=tl, pool=default, queue=default, priority_weight=7, operator=PythonOperator
```

2. Screenshots of the web browser:

a. Before login:



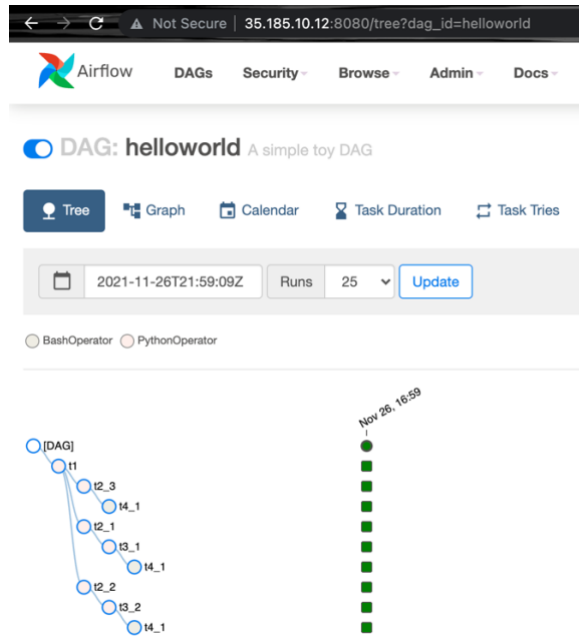
b. Showing the DAGs after successful login:



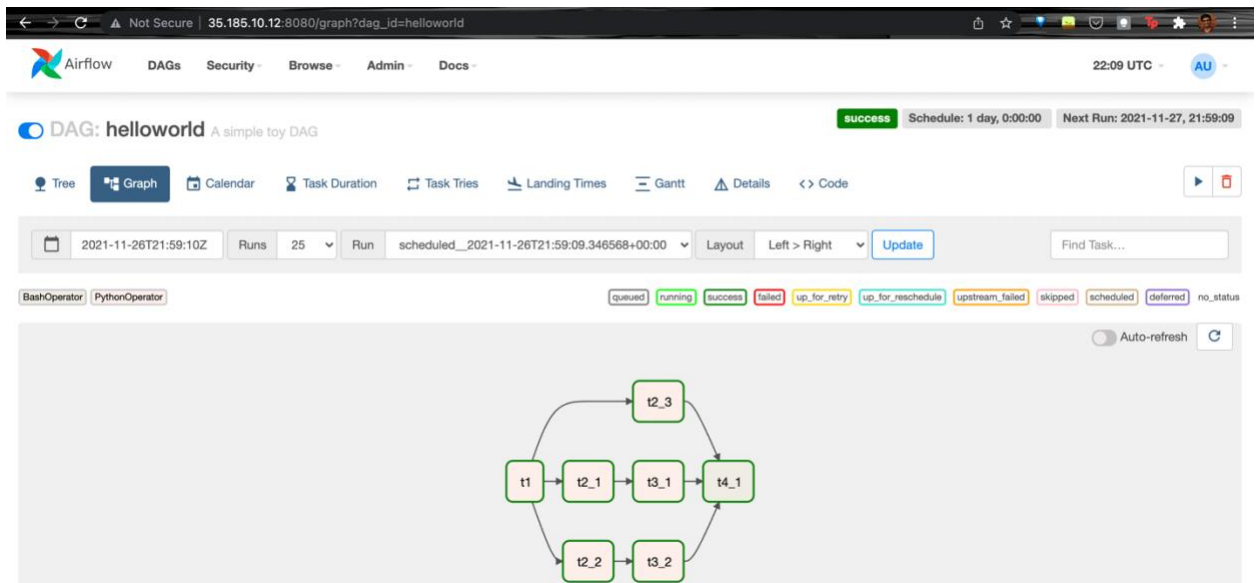
1.2 Helloworld with SequentialExecutor and LocalExecutor.

(1) SequentialExecutor:

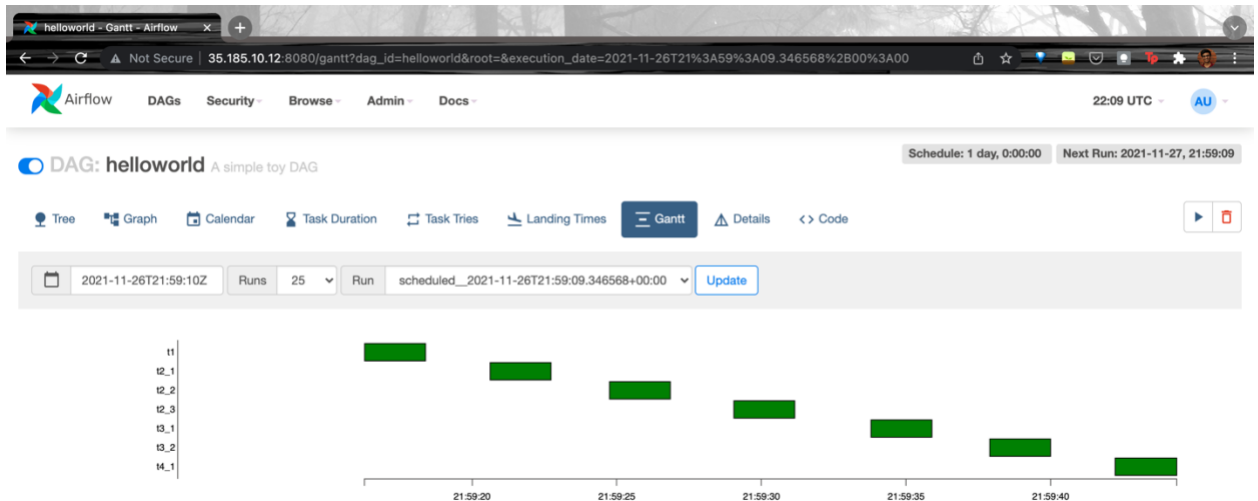
Screenshot of the tree for helloworld for Sequential Executor:



Screenshot of the graph for helloworld for Sequential Executor:

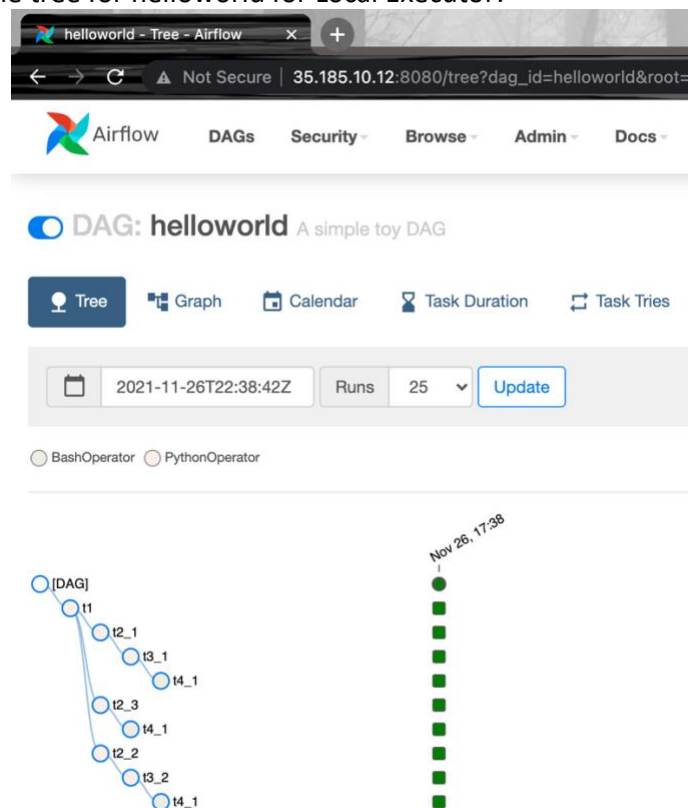


Screenshot of the gantt for helloworld for Sequential Executor:

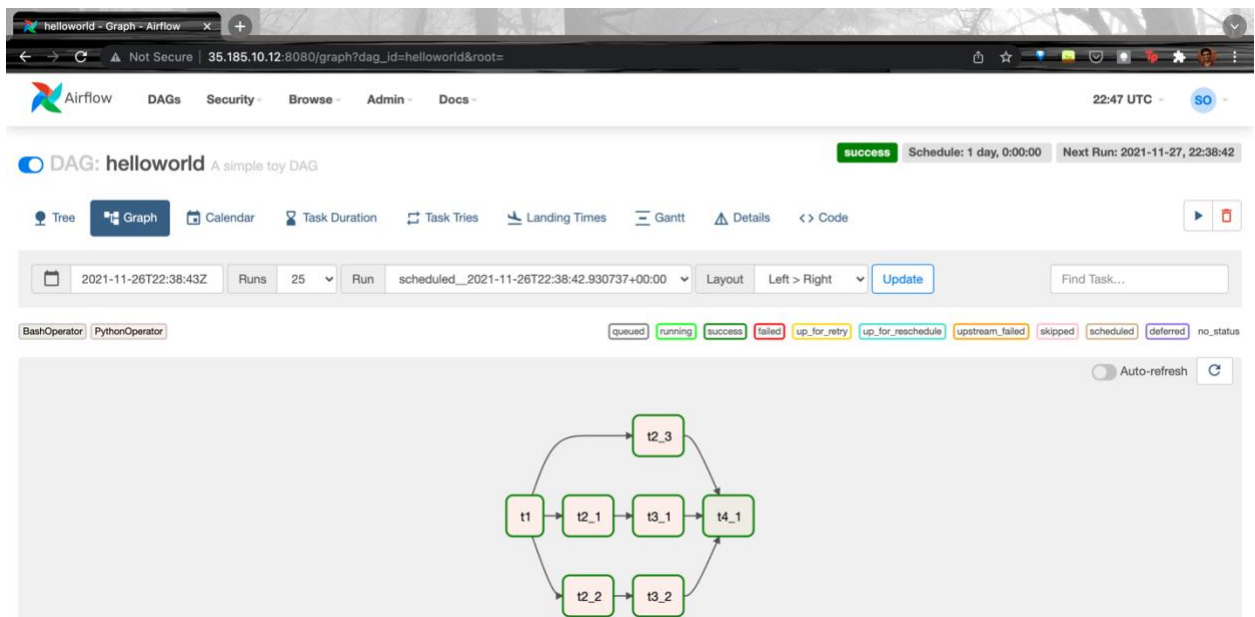


LocalExecutor:

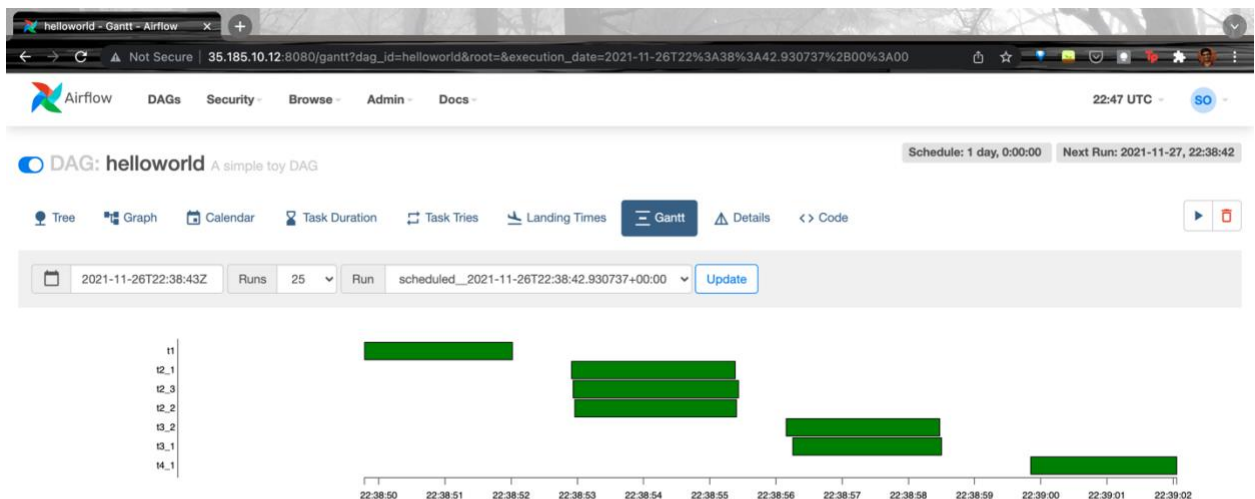
Screenshot of the tree for helloworld for Local Executor:



Screenshot of the graph for helloworld for Local Executor:



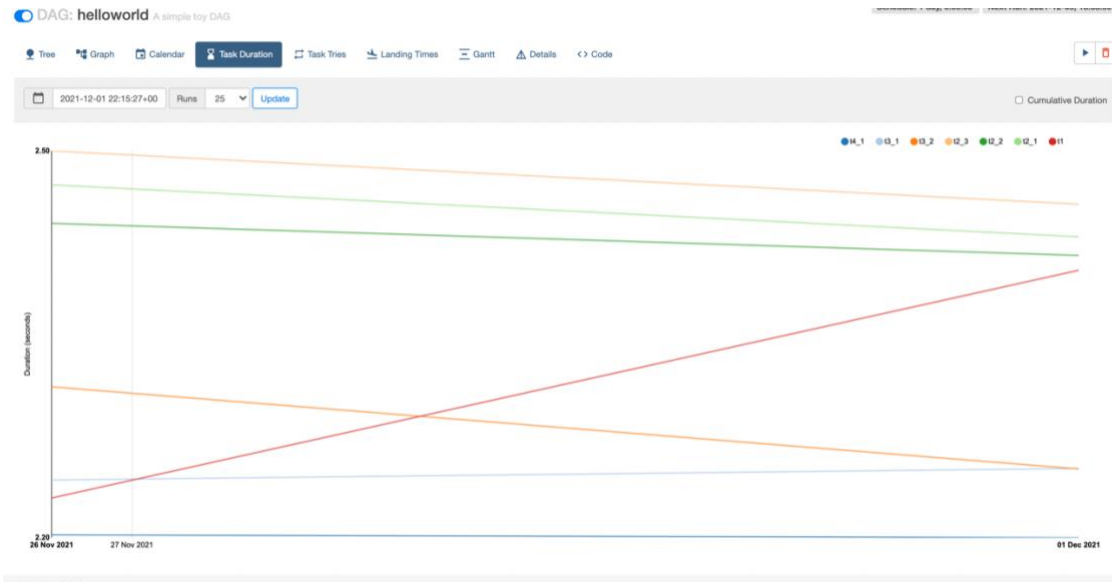
Screenshot of the gantt for helloworld for Local Executor:



(2) Two other visualizations/features in the Airflow UI are:

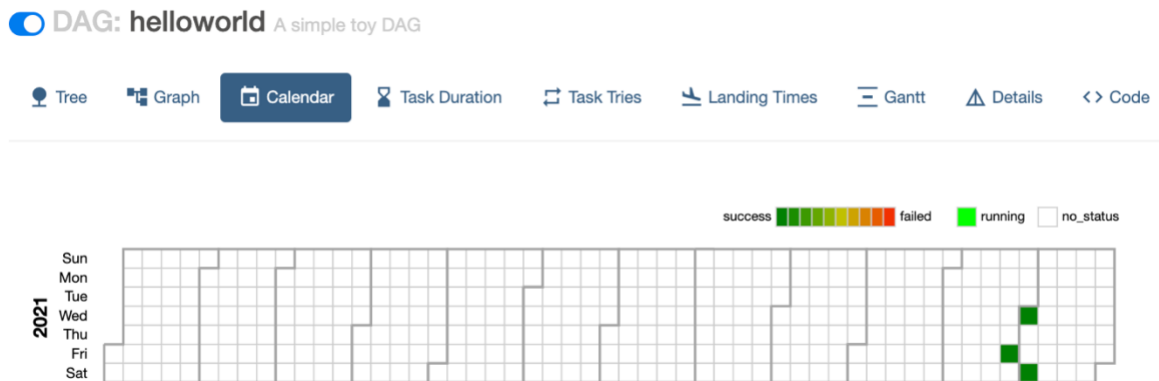
1. Task Duration: It's a visualization that shows how much time a task took to execute over the past N runs. It's a useful feature for monitoring the task execution time as it lets us find outliers and quickly understand where the time is spent in your DAG over several runs.

Screenshot of the task duration of helloworld over multiple runs is shown below:



2. Calendar: The calendar view gives us an overview of the entire history of the DAG's over months or years. It helps us in the troubleshooting process as we can understand the trends of the overall success/failure rate of runs over time from the visualisation. In the screenshot below, we can see that the helloworld DAG has tried executing three times in the past 10 days and has successfully executed all three times.

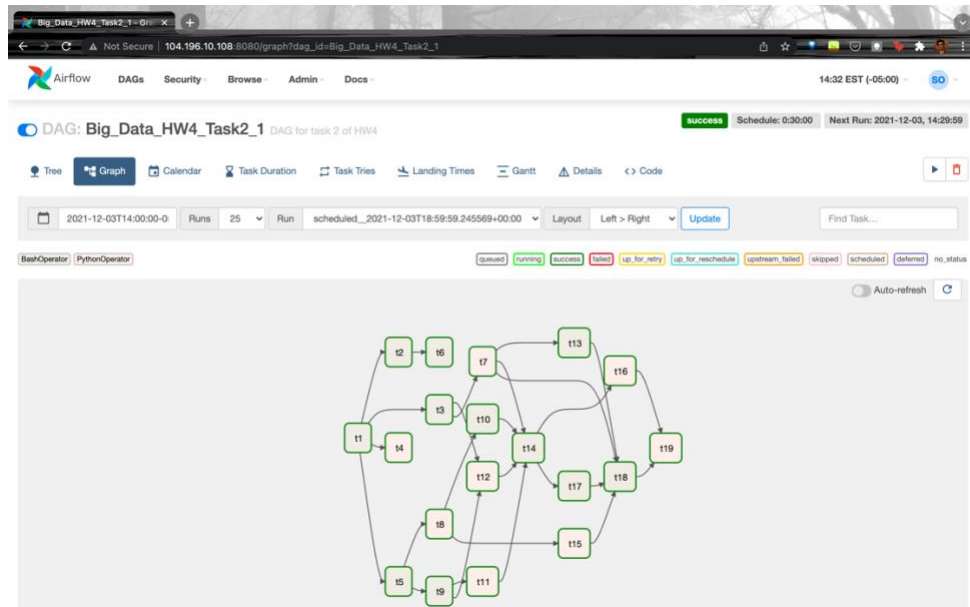
Note: The daily schedule isn't depicted here as the DAG is deployed on a virtual machine, and the virtual machine was turned off for the rest of the days.



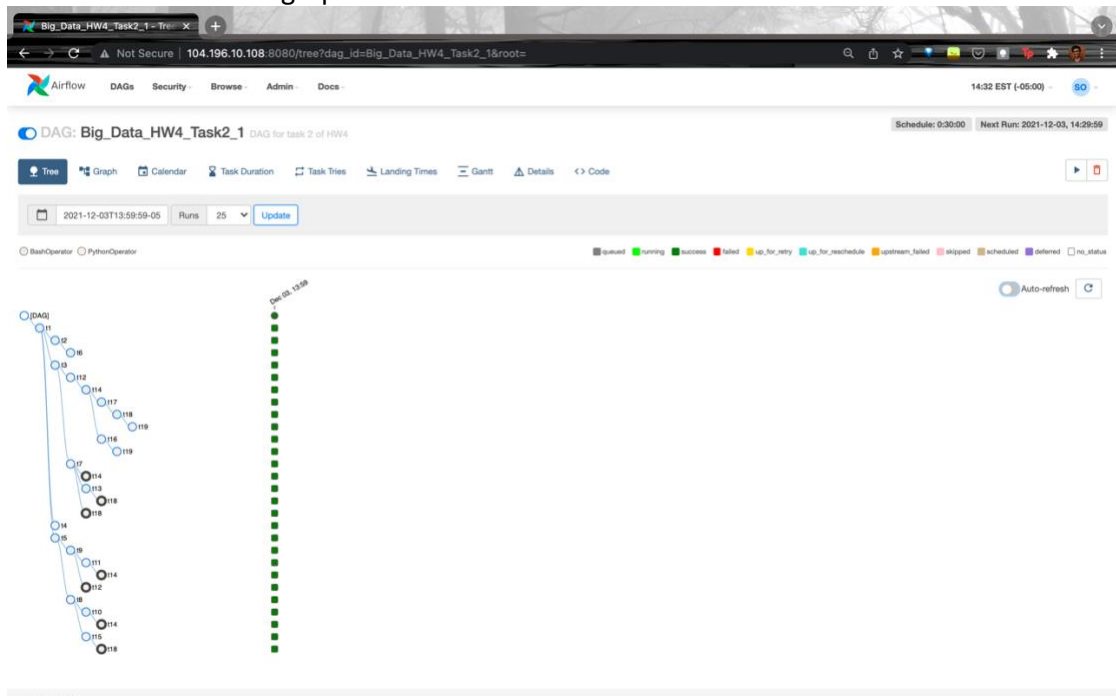
Task 2: (Build Workflows)

2.1 The DAG is implemented with a Local Executor.

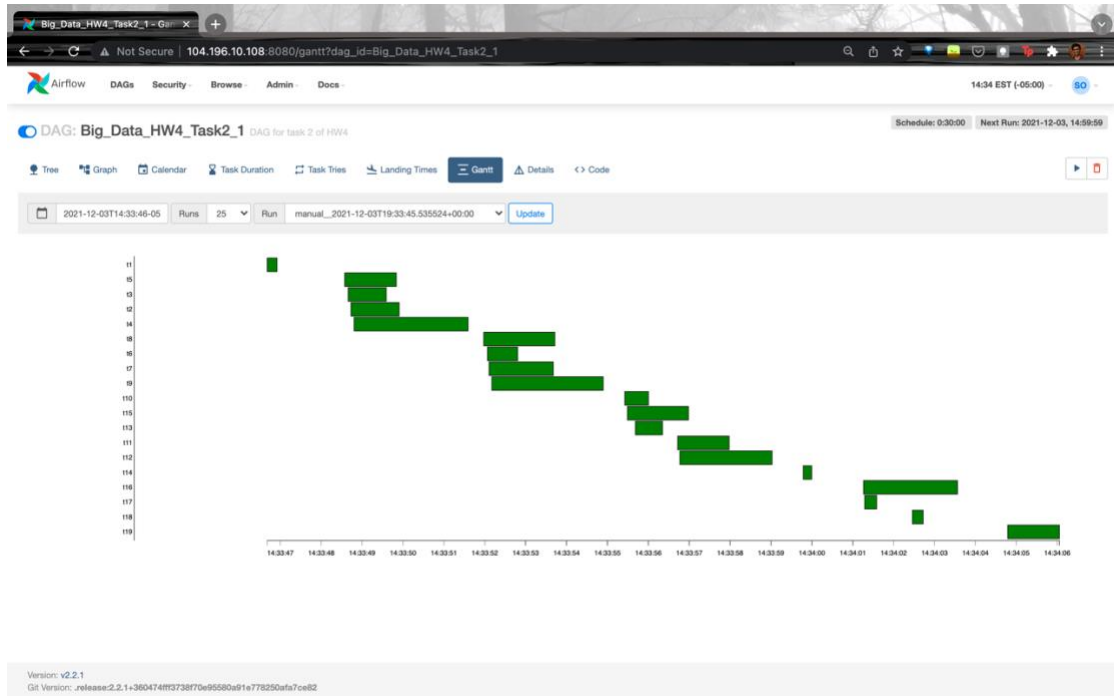
1. Screenshot of the tree:



Screenshot of the graph:



2. Screenshot pf the gantt for the manual trigger:



3. Screenshot of the running history of the DAG:

Search										
Actions										
Record Count: 5										
	State	Dag Id	Execution Date	Run Id	Run Type	Queued At	Start Date	End Date	External Trigger	Conf
<input type="checkbox"/>	success	Big_Data_HW4_Task2_1	2021-12-03, 15:30:49	scheduled_2021-12-03T20:30:49.368577+00:00	scheduled	2021-12-03, 16:00:49	2021-12-03, 16:00:49	2021-12-03, 16:01:11	False	()
<input type="checkbox"/>	success	Big_Data_HW4_Task2_1	2021-12-03, 15:00:48	scheduled_2021-12-03T20:00:48.900448+00:00	scheduled	2021-12-03, 15:30:48	2021-12-03, 15:30:48	2021-12-03, 15:31:08	False	()
<input type="checkbox"/>	success	Big_Data_HW4_Task2_1	2021-12-03, 14:33:45	manual_2021-12-03T19:33:45.535524+00:00	manual	2021-12-03, 14:33:45	2021-12-03, 14:33:46	2021-12-03, 14:34:07	True	()

As you can see above, the DAG is manually triggered at 2:33 pm EST on 3rd December. Then the next two runs are the scheduled runs, which happen every 30 minutes. The date at which the DAG starts being scheduled is called the start date. Schedule interval refers to the interval of time at which your DAG gets triggered. In the code, I have taken the start date as 2021-12-03 3:30 pm and the schedule has been set to minutes = 30, using timedelta functionality of datetime package in python. Another option to schedule would have been to use a cron job. This 30 minutes schedule interval ensures that the scheduler in the airflow triggers (manually/scheduled) the task at every 30 minutes after the set start date.

Code for Task 2.1:


```

Task2_1.py > ...
1  from datetime import datetime, timedelta
2  from textwrap import dedent
3  import time
4
5  # The DAG object; we'll need this to instantiate a DAG
6  from airflow import DAG
7
8  # Operators; we need this to operate!
9  from airflow.operators.bash import BashOperator
10 from airflow.operators.python import PythonOperator
11
12 # These args will get passed on to each operator
13 # You can override them on a per-task basis during operator initialization
14 default_args = {
15     'owner': 'Shivam',
16     'depends_on_past': False,
17     'email': ['so2639@columbia.edu'],
18     'email_on_failure': False,
19     'email_on_retry': False,
20     'retries': 1,
21     'retry_delay': timedelta(minutes=30)
22 }
23
24
25 count = 0
26
27 def sleeping_function():
28     """This is a function that will run within the DAG execution"""
29     time.sleep(2)
30
31 def count_function():
32     global count
33     count += 1
34     print('count_increase output: {}'.format(count))
35     time.sleep(1)
36
37 def print_function():
38     print('This task is a python operator')
39     time.sleep(1)
40
41

```

```

Task2_1.py > ...
41
42 with DAG(
43     'Big_Data_HW4_Task2_1',
44     default_args=default_args,
45     description='DAG for task 2 of HW4',
46     schedule_interval=timedelta(minutes=30),
47     start_date=datetime(2021, 12, 3, 3, 30, 00),
48     catchup=False,
49     tags=['Task2'],
50 ) as dag:
51
52     # task examples of tasks created by instantiating operators
53
54     t1 = BashOperator(
55         task_id='t1',
56         bash_command='echo "This task is part of Question 2"',
57     )
58
59     t2 = BashOperator(
60         task_id='t2',
61         bash_command='python3 /home/so2639/airflow/dags/py_script.py',
62         retries=2,
63     )
64
65     t3 = BashOperator(
66         task_id='t3',
67         bash_command='date',
68         retries=3,
69     )
70
71     t4 = PythonOperator(
72         task_id='t4',
73         python_callable=sleeping_function,
74     )
75
76     t5 = BashOperator(
77         task_id='t5',
78         bash_command='python3 /home/so2639/airflow/dags/py_script.py',
79     )
80
81     t6 = BashOperator(

```

```

Task2_1.py > ...
80
81     t6 = BashOperator(
82         task_id='t6',
83         bash_command='echo "This task is part of Question 2"',
84     )
85
86     t7 = PythonOperator(
87         task_id='t7',
88         python_callable=print_function,
89         retries=2,
90     )
91
92     t8 = PythonOperator(
93         task_id='t8',
94         python_callable=count_function,
95         retries=2,
96     )
97
98     t9 = BashOperator(
99         task_id='t9',
100        bash_command='sleep 2',
101    )
102
103    t10 = BashOperator(
104        task_id='t10',
105        bash_command='date',
106    )
107
108    t11 = PythonOperator(
109        task_id='t11',
110        python_callable=print_function,
111        retries=2,
112    )
113
114    t12 = PythonOperator(
115        task_id='t12',
116        python_callable=sleeping_function,
117        retries=2,
118    )
119
120    t13 = BashOperator(

```

```

Task2_1.py > ...
120    t13 = BashOperator(
121        task_id='t13',
122        bash_command='date',
123    )
124
125    t14 = BashOperator(
126        task_id='t14',
127        bash_command='echo "This task is part of Question 2"',
128    )
129
130    t15 = PythonOperator(
131        task_id='t15',
132        python_callable=count_function,
133        retries=2,
134    )
135
136    t16 = PythonOperator(
137        task_id='t16',
138        python_callable=sleeping_function,
139        retries=2,
140    )
141
142    t17 = BashOperator(
143        task_id='t17',
144        bash_command='echo "This task 17 is part of Question 2"',
145    )
146
147    t18 = BashOperator(
148        task_id='t18',
149        bash_command='python3 /home/so2639/airflow/dags/py_script.py',
150    )
151
152    t19 = PythonOperator(
153        task_id='t19',
154        python_callable=print_function,
155        retries=1,
156    )
157
158    # task dependencies
159    t1 >> [t2, t3, t4, t5]
160

```

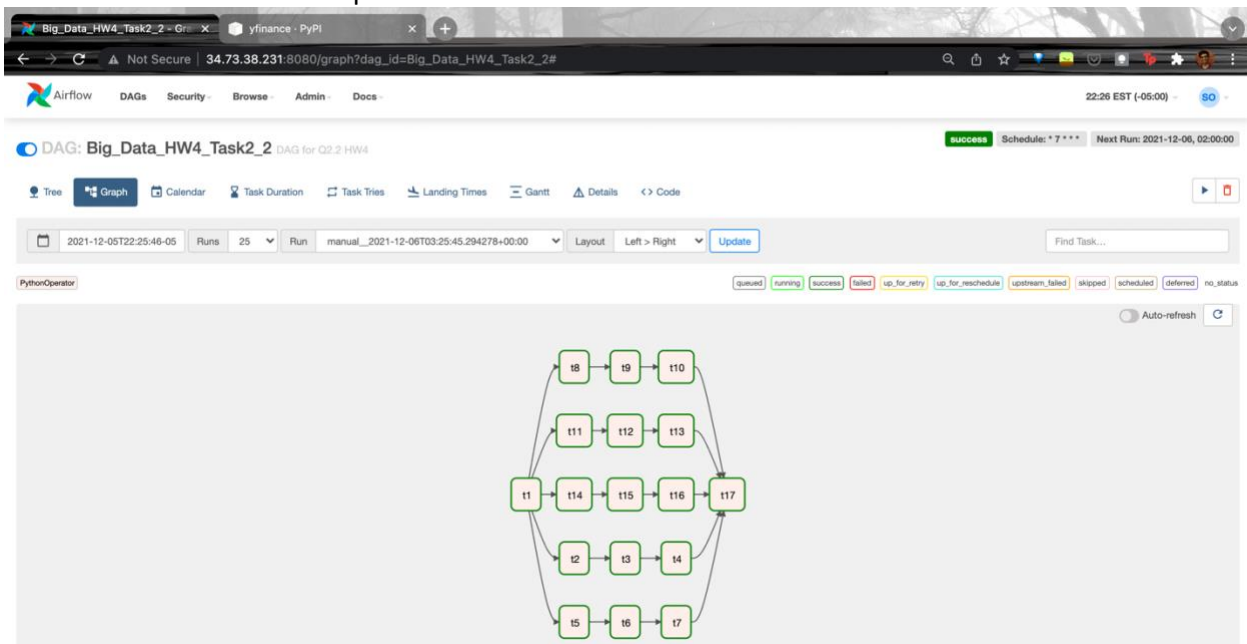
```

Task2_1.py > ...
151
152     t19 = PythonOperator(
153         task_id='t19',
154         python_callable=print_function,
155         retries=1,
156     )
157
158
159     # task dependencies
160     t1 >> [t2, t3, t4, t5]
161     t2 >> t6
162     t3 >> [t7, t12]
163     t5 >> [t8, t9]
164     t7 >> [t13, t14, t18]
165     t8 >> [t10, t15]
166     t9 >> [t11, t12]
167     [t10, t11, t12] >> t14
168     [t13, t15, t17] >> t18
169     t14 >> [t16, t17]
170     [t16, t18] >> t19
171
172

```

2.2 Stock price fetching, prediction and storage:

Airflow DAG Graph View:



1. I have used the yfinance python package to fetch the data and created a DAG on airflow and set a cron job '* * 7 * * *' to schedule it to run at 7:00 AM every day (shown above).
2. The data is pre-processed by dropping the nan values, if any, with the help of dropna method in pandas, so that the model training later on isn't affected.

- As per the question, 5 linear regression models are trained(each for one company) with the the features as 'high price', 'low price', 'close price', open price', 'volume'] and it predicts the 'high price' for the next day.
- An if condition checks if the values received are same for the last 2 days(They can be same in case the market was closed). If they are different, the relative error for each stock is predicted and stored in a output csv. A final compiled csv is created at the end that shows all the error values for all the companies for each day the DAG has been running. The relative error is calculated from 30th November to 5th December for each company, shown in the csv table screenshot below:

	A	B	C	D	E	F	G	H	I	J
1		Date	AAPL	GOOGL	FB	MSFT	AMZN			
2	0	11/30/21	0.01833223	0.0178464	0.02154729	0.006656	-0.0011046			
3	0	12/1/21	-0.018877	-0.0186309	-0.0090627	-0.0187762	-0.0092226			
4	0	12/2/21	-0.0229995	0.00372933	0.02005463	0.00508727	0.00354409			
5	0	12/3/21	-0.0196708	-0.0153219	-0.0053576	-0.0123388	-0.0057598			
6	0	12/4/21	0.01860972	-0.0068104	0.00302498	0.00552341	-0.0008707			
7	0	12/5/21	0.01860972	-0.0068104	0.00302498	0.00552342	-0.0008707			
8										

- The workflow is setup using a DAG. Each functionality of the code is defined as a function, and different tasks in the DAG call these functions one by one. The stock data for the five companies is fetched and processed in parallel (Task 2, 5, 8, 11 and 14 start in parallel, and the subsequent pre-processing and prediction tasks for each company also happen in parallel accordingly). Once the model is compiled and the predictions are ready, the final output csv is generated with the relative errors for all the companies(Task 17). The task dependencies are mentioned below:

```

t1 >> [t2, t5, t8, t11, t14]
t2 >> t3
t3 >> t4
t5 >> t6
t6 >> t7
t8 >> t9
t9 >> t10
t11 >> t12
t12 >> t13
t14 >> t15
t15 >> t16
[t4, t7, t10, t13, t16] >> t17

```

For the Linear regression model training, the training data for the model is taken till two days before the day for which we have to predict the stock price. For example, consider we fetch data on 30th November 7 AM, hence we will fetch the stock prices till 29th November, when the market closed. So as per the question, we will use the data till 28th November to train the model, and as we have to make predictions for the high price for the next day (29th November), we will push the index of predictions by 1. Hence, we will be able to calculate the relative error for the last record i.e. 29th November.

The csv and json files are used for cross task communication. For each company the data is fetched, pre-processed and stored as a csv initially. This data is used to calculate the relative error in a later task, and the error is stored in a dict initially. Once the dict has accumulated the error values for all the companies, then the last task(t17) calls the generate_output_csv function, which stores the error values as a dataframe for the given date. Then all the date error values are compiled to create a single csv file 'compiled_error_results.csv'. The scheduler is setup with the help of cron job, so that the DAG executes every day at 7 am and fetches the stock data till previous day (As at 7 am, the stock price for the current day will not be available).

Screenshots for code for Task 2.2:

```
so2639_6893_hw4_Task2_2.py 6 X
so2639_6893_hw4_Task2_2.py > preprocess_data
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  import os
5  import json
6  import yfinance as yf
7  import numpy as np
8  import pandas as pd
9  from datetime import datetime, timedelta
10 from airflow import DAG
11 from airflow.operators.python import PythonOperator
12 from sklearn.linear_model import LinearRegression
13
14 default_args = {
15     'owner': 'Shivam',
16     'depends_on_past': False,
17     'email': ['so2639@columbia.edu'],
18     'email_on_failure': False,
19     'email_on_retry': False,
20     'retries': 1,
21     'retry_delay': timedelta(seconds=30)
22 }
23
24 data_dir = '/home/so2639/airflow/dags/task2_2/'
25 json_file = 'errors_dict.json'
26 date = datetime.now().strftime("%Y-%m-%d")
27 #strt_date = datetime.strptime(date, "%Y-%m-%d") - timedelta(weeks=27)
28
29 def fetch_data(ticker):
30     input_company_data = yf.Ticker(ticker)
31     # Fetching data for the last 6 months
32     # df = input_company_data.history(start=strt_date, end=date)
33     df = input_company_data.history(period = '6mo')
34     df.to_csv(data_dir + ticker + '.csv')
35
36 def preprocess_data(ticker):
37     df = pd.read_csv(data_dir + ticker + '.csv')
38     df['Date'] = pd.to_datetime(df['Date'])
39     # The data doesn't necessarily need to be normalised
40     # Cleaning the data by dropping null values and any date values greater than today
41     df = df[df.Date <= pd.to_datetime(date)].reset_index(drop = True)
42     df = df.dropna()
43     df.to_csv(data_dir + ticker + '.csv')
```

so2639_6893_hw4_Task2_2.py 6 X

so2639_6893_hw4_Task2_2.py > preprocess_data

```
43 | df.to_csv(data_dir + ticker + '.csv')
44 |
45 | def init_df_dict():
46 |     errors_dict = {'Date': '', 'AAPL': '', 'GOOGL': '', 'FB': '', 'MSFT': '', 'AMZN': ''}
47 |     errors_dict['Date'] = [date]
48 |
49 |     with open(data_dir + json_file, 'w') as jsonfile:
50 |         json.dump(errors_dict, jsonfile)
51 |
52 | def predict_and_compute_error(ticker):
53 |     df = pd.read_csv(data_dir + ticker + '.csv')
54 |     df = df.iloc[:-1, :]
55 |
56 |     # Generate training data
57 |     X = np.array(df[['Open', 'High', 'Low', 'Close', 'Volume']])
58 |     y = np.array(df['High'])
59 |
60 |     # When we fetch data till 30th November morning, we get stock values till market close on 29th.
61 |     # So using data till 28th November to train. Pushing the labels value ahead by 1 index as well.
62 |     # error = (prediction yesterday - actual price today)/actual price today
63 |     # example: error = (prediction on 28th for 29th - price on 29th)/ price on 29th
64 |     X_train = X[:-2, :]
65 |     y_train = y[1:-1]
66 |
67 |     # Run Linear Regression Model
68 |     lr_model = LinearRegression().fit(X_train, y_train)
69 |     print("Score: " + str(lr_model.score(X_train, y_train)))
70 |
71 |     # Make predictions
72 |     X_test = X[-2:-1, :] # Last value in dataframe
73 |     y_test = y[-1] # Last High Value
74 |
75 |     preds = lr_model.predict(np.array(X_test))
76 |     preds = preds[0]
77 |
78 |     # Read error json
79 |     with open(data_dir + json_file) as j:
80 |         errors = json.load(j)
81 |
```


so2639_6893_hw4_Task2_2.py 6 ×

so2639_6893_hw4_Task2_2.py > preprocess_data

```
82     # Compute error
83     if not(np.array_equal(X[-2:-1, :], X[-3:-2, :], equal_nan=False)):
84         er = (preds - y_test) / y_test
85         errors[ticker] = [er]
86
87     # Update error json
88     with open(data_dir + json_file, 'w') as j:
89         json.dump(errors, j)
90
91 def generate_output_csv():
92     # Read error json
93     with open(data_dir + json_file) as j:
94         errors = json.load(j)
95
96     df = pd.DataFrame(errors)
97     df.to_csv(data_dir + 'csv_files/' + 'results_' + date + '.csv')
98
99     csv_list = []
100    for file in os.listdir(data_dir+'csv_files/'):
101        if (file.startswith('results_') and file.endswith('.csv')):
102            csv_list.append(os.path.join(data_dir + 'csv_files/', file))
103
104    # Compute compiled csv output
105    out = pd.DataFrame(columns=['Date', 'AAPL', 'GOOGL', 'FB', 'MSFT', 'AMZN'])
106    for item in csv_list:
107        df1 = pd.read_csv(item, index_col=[0])
108        out = out.append(df1)
109    out.sort_values(by=['Date'], inplace=True, ascending=True)
110    out.to_csv(data_dir + 'compiled_error_results.csv')
111
112    with DAG(
113        'Big_Data_HW4_Task2_2',
114        default_args = default_args,
115        description = 'DAG for Q2.2 HW4',
116        start_date = datetime(2021, 11, 30, 7, 0, 0),
117        schedule_interval = '* 7 * * *',
118        catchup = False,
119        tags = ['example'],
120    ) as dag:
121
```


so2639_6893_hw4_Task2_2.py 6 X

so2639_6893_hw4_Task2_2.py > ...

```
121
122     t1 = PythonOperator(
123         task_id = 't1',
124         python_callable = init_df_dict,
125     )
126
127     t2 = PythonOperator(
128         task_id = 't2',
129         python_callable = fetch_data,
130         op_kwargs = {'ticker': 'AAPL'}
131     )
132
133     t3 = PythonOperator(
134         task_id = 't3',
135         python_callable = preprocess_data,
136         op_kwargs = {'ticker': 'AAPL'}
137     )
138
139     t4 = PythonOperator(
140         task_id = 't4',
141         python_callable = predict_and_compute_error,
142         op_kwargs = {'ticker': 'AAPL'}
143     )
144
145     t5 = PythonOperator(
146         task_id = 't5',
147         python_callable = fetch_data,
148         op_kwargs = {'ticker': 'GOOGL'}
149     )
150
151     t6 = PythonOperator(
152         task_id = 't6',
153         python_callable = preprocess_data,
154         op_kwargs = {'ticker': 'GOOGL'}
155     )
156
157     t7 = PythonOperator(
158         task_id = 't7',
159         python_callable = predict_and_compute_error,
160         op_kwargs = {'ticker': 'GOOGL'}
161     )
```

so2639_6893_hw4_Task2_2.py 6 ×

so2639_6893_hw4_Task2_2.py > ...

```
161 )
162
163 t8 = PythonOperator(
164     task_id = 't8',
165     python_callable = fetch_data,
166     op_kwargs = {'ticker': 'FB'}
167 )
168
169 t9 = PythonOperator(
170     task_id = 't9',
171     python_callable = preprocess_data,
172     op_kwargs = {'ticker': 'FB'}
173 )
174
175 t10 = PythonOperator(
176     task_id = 't10',
177     python_callable = predict_and_compute_error,
178     op_kwargs = {'ticker': 'FB'}
179 )
180
181
182 t11 = PythonOperator(
183     task_id = 't11',
184     python_callable = fetch_data,
185     op_kwargs = {'ticker': 'MSFT'}
186 )
187
188 t12 = PythonOperator(
189     task_id = 't12',
190     python_callable = preprocess_data,
191     op_kwargs = {'ticker': 'MSFT'}
192 )
193
194 t13 = PythonOperator(
195     task_id = 't13',
196     python_callable = predict_and_compute_error,
197     op_kwargs = {'ticker': 'MSFT'}
198 )
199
```

```

so2639_6893_hw4_Task2_2.py 6 ×
so2639_6893_hw4_Task2_2.py > ...
200
201     t14 = PythonOperator(
202         task_id = 't14',
203         python_callable = fetch_data,
204         op_kwargs = {'ticker': 'AMZN'}
205     )
206
207     t15 = PythonOperator(
208         task_id = 't15',
209         python_callable = preprocess_data,
210         op_kwargs = {'ticker': 'AMZN'}
211     )
212
213     t16 = PythonOperator(
214         task_id = 't16',
215         python_callable = predict_and_compute_error,
216         op_kwargs = {'ticker': 'AMZN'}
217     )
218
219     t17 = PythonOperator(
220         task_id = 't17',
221         python_callable = generate_output_csv,
222     )
223
224     # Task dependencies
225     t1 >> [t2, t5, t8, t11, t14]
226     t2 >> t3
227     t3 >> t4
228     t5 >> t6
229     t6 >> t7
230     t8 >> t9
231     t9 >> t10
232     t11 >> t12
233     t12 >> t13
234     t14 >> t15
235     t15 >> t16
236     [t4, t7, t10, t13, t16] >> t17
237

```

Task 3: (Written Parts)

3.1 Sequential Executor:

Pros: Sequential Executor allows to run Airflow without setting up many dependencies. So, we can even work directly with SQLite using the SequentialExecutor. It also identifies a single point of failure and is hence easy for debugging.

Cons: It is not recommended for any use cases that require more than a single task execution at a time, and hence usually it cannot be used in production.

Local Executor:

Pros: Local Executor is easy to set up and is resource efficient. It allows testing multiple jobs in parallel & hence offers parallelism.

Cons: It is not sufficiently scalable, and hence is only used in small-scale production environments. It is also affected by a single point of failure.

Celery Executor:

Pros: Celery Executor quickly adapts and is able to assign that allocated task or tasks to another worker, if a worker node is ever down or goes offline. It is built for horizontal scaling and has high availability.

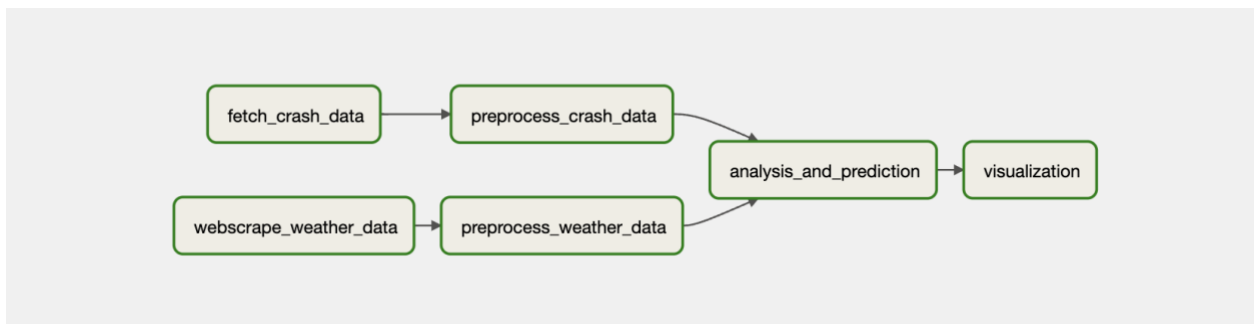
Cons: It is expensive and worker maintenance is required. It's set up is also comparatively complicated.

Kubernetes Executor:

Pros: Kubernetes Executor is highly resource & cost efficient. Here, each task instance is run in its own pod on a Kubernetes cluster. It runs as a process in the Airflow Scheduler. When the traffic is high, we can scale up. Conversely, when the traffic is low, we can scale to zero. It also has high level of fault tolerance and has task-level configurations. Also, if a deployment is pushed, there are no interruptions to running tasks.

Cons: There is an overhead of a few extra seconds per task for a pod to spin up. Also lack of Kubernetes familiarity can potentially act as a barrier to use this.

3.2 DAG of group project “NYC Vehicle Crash Analysis”:



1) The DAG for our group project “NYC Vehicle Crash Analysis” is sequenced into 6 tasks which is shown in the graph from the airflow UI above. As depicted, the tasks have been divided based on the various steps in our ETL pipeline and the analysis and visualization part of our project.

2) Tasks:

- a. fetch_crash_data: The complete historical and incremental data load of the crash dataset is done using the SODA API available on the NYC Open Data website. This formulates the first task of the DAG, which runs the python script to do the operations mentioned above.
- b. fetch_weather_data: This task executes the web scraping script to fetch the weather data. The fetched raw data is stored in the Cloud SQL postgres server

and this task executes in parallel to the crash dataset to reduce the total execution time of the pipeline.

- c. preprocess_crash_data: This step involves taking the raw crash data and backfilling the missing zip code and latitude/longitude values.
- d. preprocess_weather_data: This step involves taking the raw weather data and cleaning up the relevant columns that will be used in the recurrent neural network later on to predict the future weather conditions.
- e. analysis_and_prediction: This step involves making correlations of crash and weather data and making weather predictions for the next two months. These predictions are used along with the crash data to predict the chances of crashes happening in a borough in New York city in the next two months.
- f. visualization: The analysis done in the previous step is visualized using metabase in this task. The data used here so be fully processed so that we are able to make a better analysis. Hence this is the last task of our project.

Task Dependencies:

The dependencies for the DAG are depicted below:

```
fetch_crash_data >> preprocess_crash_data
fetch_weather_data >> preprocess_weather_data
[preprocess_crash_data, preprocess_weather_data] >> analysis_and_prediction
analysis_and_prediction >> visualization
```

As mentioned earlier, the tasks have been divided based on the several steps in the ETL pipeline and the analysis and visualization part of our project. Analysis and visualization tasks have upstream dependencies, as mentioned above, to make sure that the data used for the final analysis and visualization is clean and processed.

- 3) The airflow DAG is scheduled to run daily at 7pm, so that it is able to fetch the incremental crash and weather data and update the analysis on a daily basis. This schedule is set in the airflow arguments in the code, where the schedule_interval is mentioned as days = 1, and start time is at 15th November 2021 at 7 pm.