

# Final Project Assignment

## DADS 6005 : Data Streaming and Real-Time Analytics

Chayaphon Sornthananon, 6610422007

Semester 1/2024

Lecturer : Asst.Prof Ekarat Rattagan



# Table of Content

<u>Requirements</u>	P.3
<u>Airflow DAG Diagram</u>	P.4
<u>DAG 1 : Get BTC and Save to MongoDB</u>	P.5-6
<u>DAG 2 : Train Model</u>	P.7-9
<u>DAG 3 : Predict BTC Price</u>	P.10-12
<u>Appendix</u>	P.13-20

# Requirements

รันระบบ airflow สามวัน ดังนี้ วันที่ 21, 22, 23 ธค โดยในแต่ละวันจะกำหนดช่วงเวลาในการ predict BTC ตั้งแต่ 22:00 ~ 6:00 โมงเช้า โดย schedule ทุกๆ 5 นาที และคำนวณ mape ทุกๆ ชั่วโมง (หมายถึง len=12) โดย นศ จะส่ง mape.txt 3 files (ของช่วงวันเวลาที่ทดสอบ 3 วัน) + code 3 dags + requirement.txt และ files อื่นๆ ถ้ามีความแตกต่างจาก starter kit กำหนด

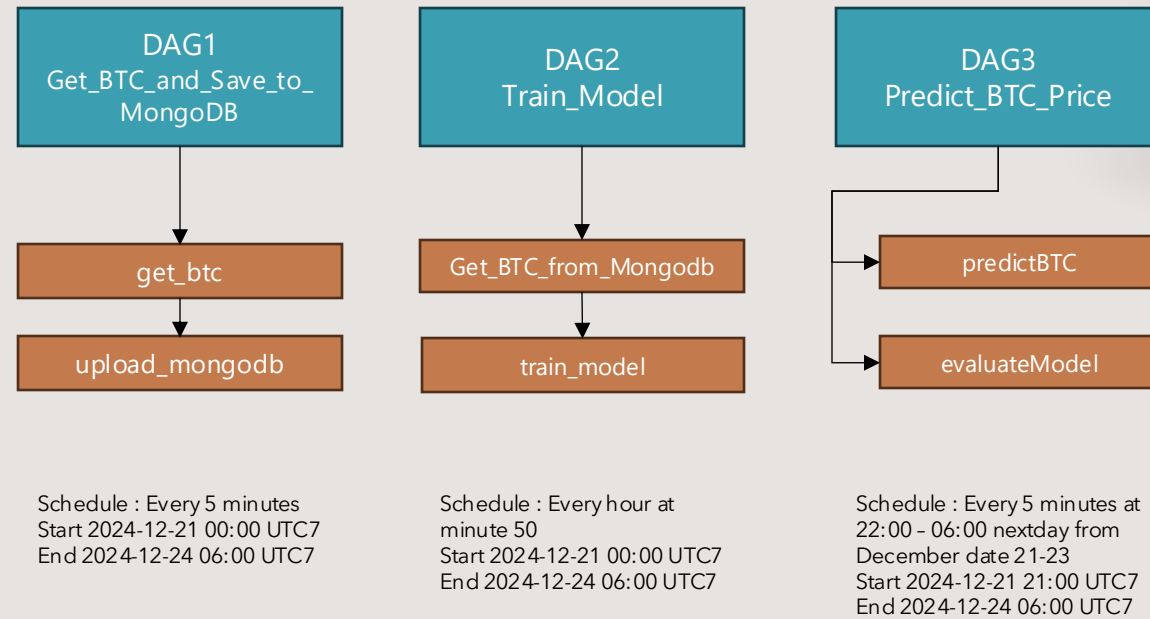
1. dataset ในการสร้าง train และ predict model จะมี column (X=closeTime) และ (Y=lastPrice)
2. นศ สามารถใช้ regression model ต่างๆ หรือ เทคนิคทางด้าน time series หรือ AI มาช่วยวิเคราะห์ได้

# Airflow DAG Diagram

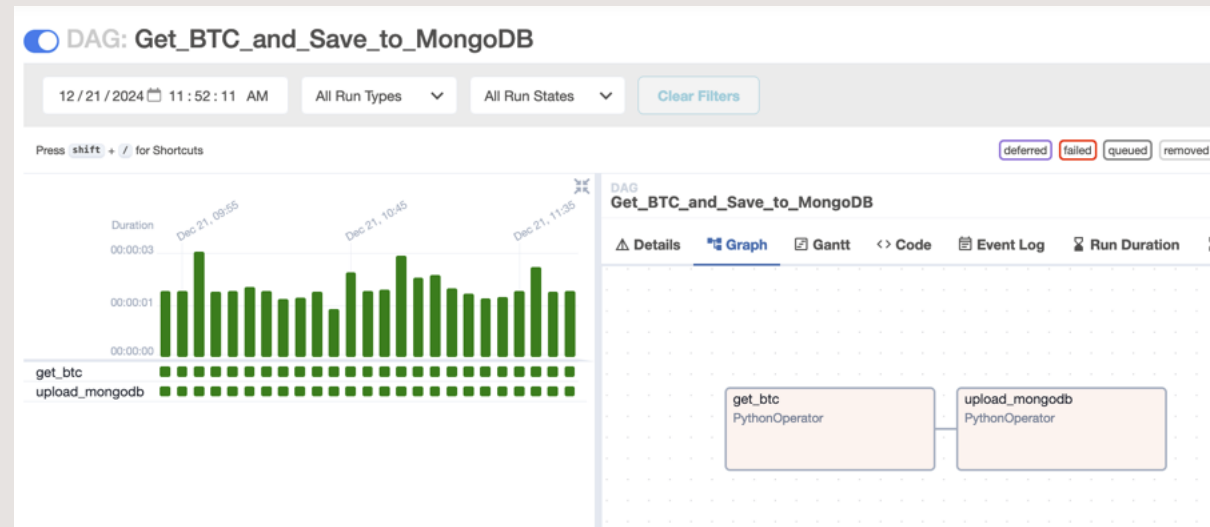
DAG

Task

Schedule



# DAG 1 : Get\_BTC\_and\_Save\_to\_MongoDB



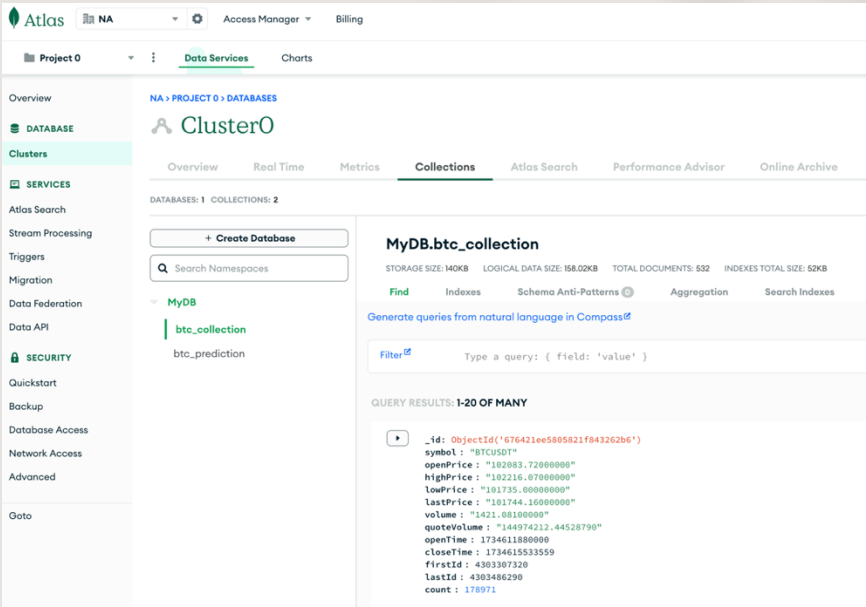
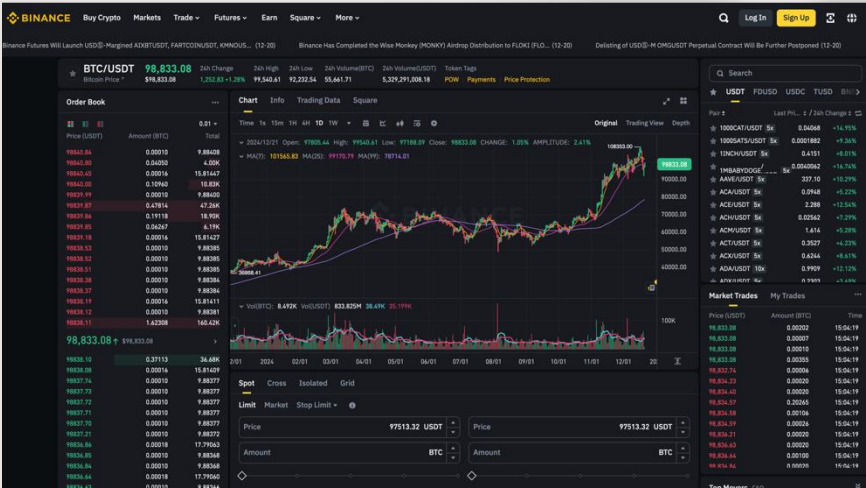
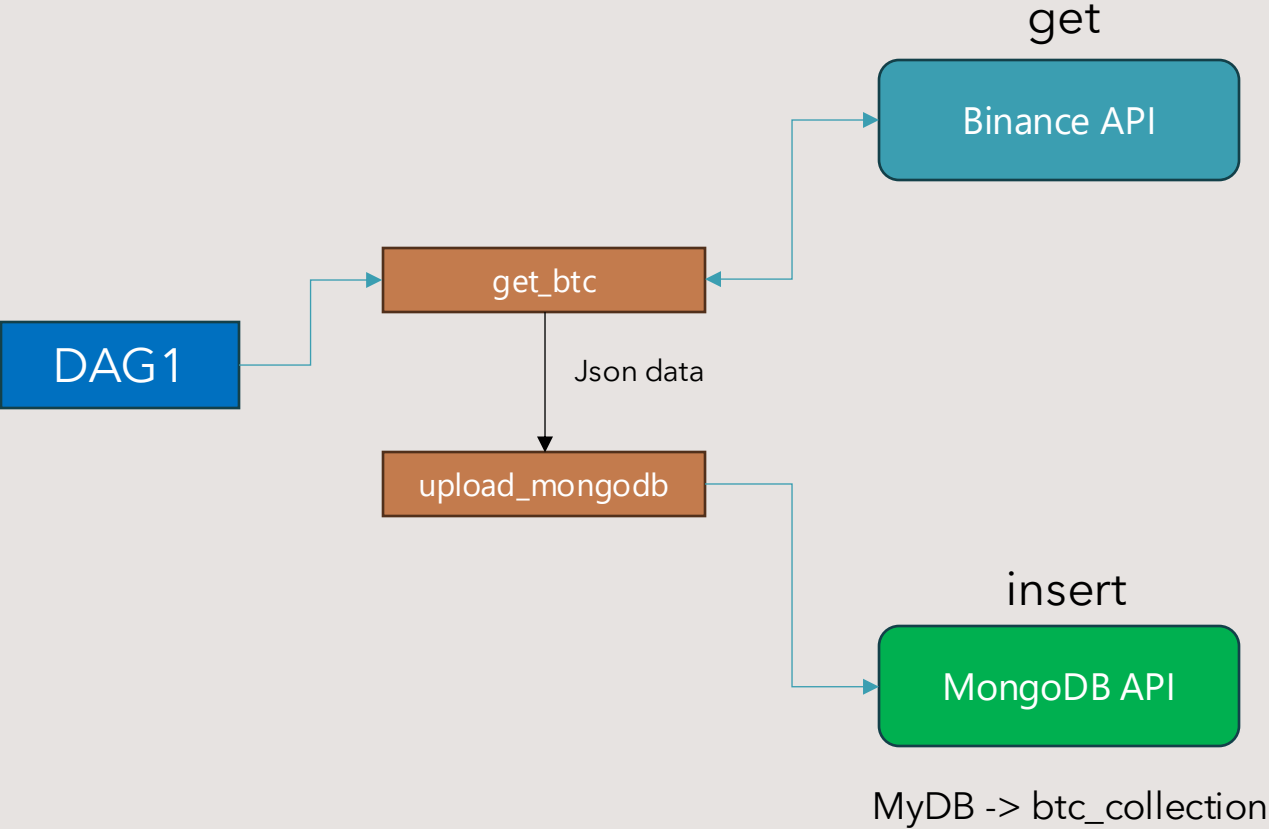
There are 2 task in DAG 1 as following :

- **Get\_btc** : This task use [https requests](#) to get current BTCUSDT price from Binance API with windows size of 1 hours (average last 1 hour metrics) which contain data ('\_id', 'symbol', 'openPrice', 'highPrice', 'lowPrice', 'lastPrice', 'volume', 'quoteVolume', 'openTime', 'closeTime', 'firstId', 'lastId', 'count') in json format.
- **Upload\_mongodb** : This task insert the json data output from above task to MongoDB cloud (MyDB->btc\_collection)

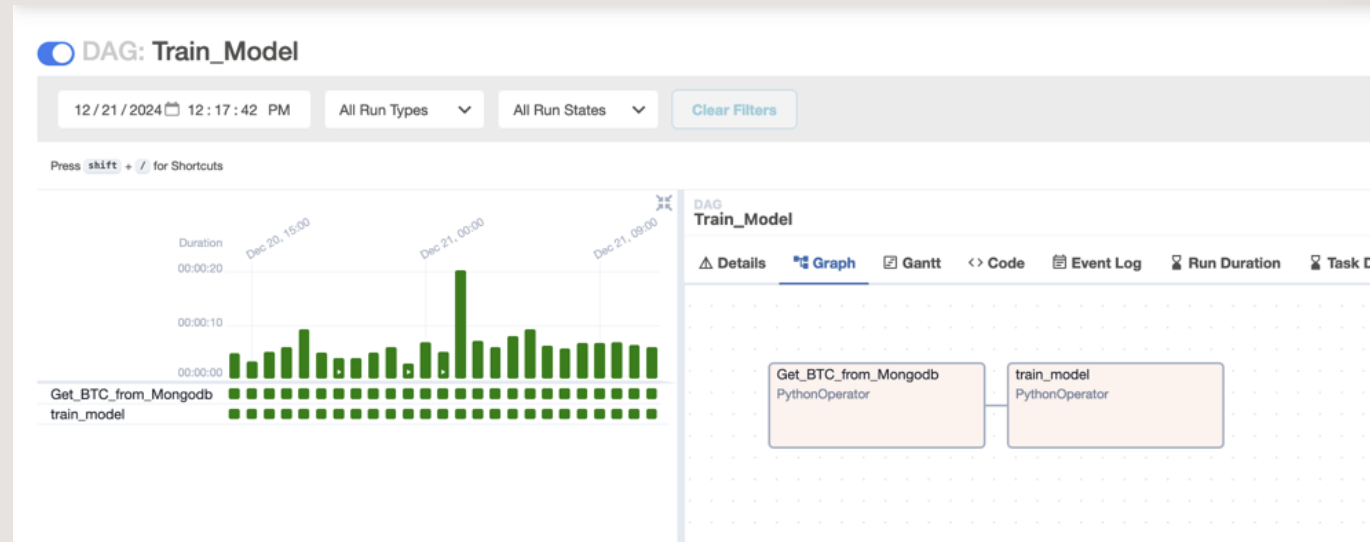
Schedule : Every 5 minutes. start 2024-12-21 00:00 utc7, end 2024-12-24 06:00 utc7

*\*More detail please check on code at dag1.py*

# DAG 1 : Get\_BTC\_and\_Save\_to\_MongoDB



## DAG 2 : Train\_Model



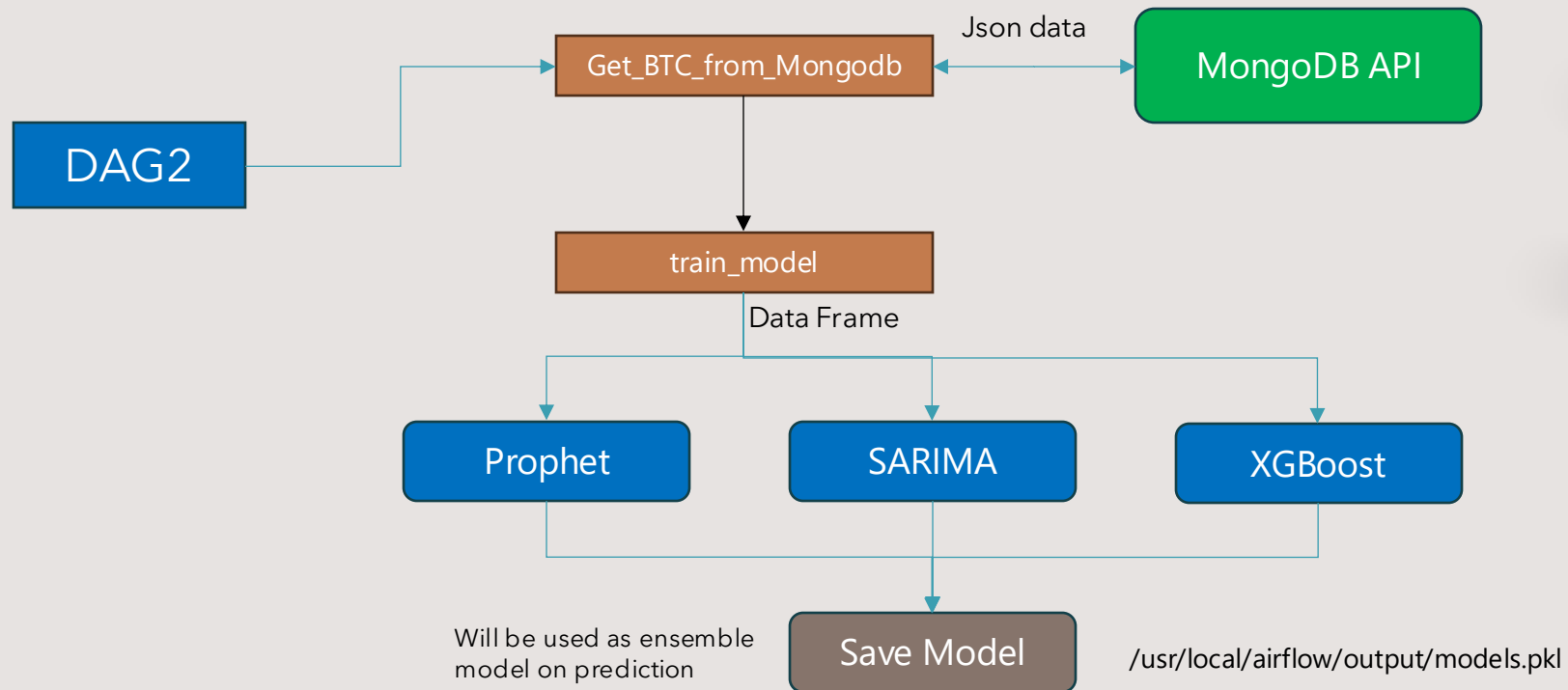
There are 2 task in DAG 2 as following :

- **Get\_BTC\_from\_Mongodb** : This task load data (last1day from current) from MongoDB cloud (MyDB->btc\_collection).
- **Train\_model** : use data output from above task (last 1day from current) with closeTime(X) and lastPrice(y) to train model.  
Then save model to “/usr/local/airflow/output/models.pkl”  
This project come with the idea of **ensemble model** of (Prophet, SARIMA and XGBoost) so we train all these 3 models in once.
  - **Prophet**: Suitable for handling seasonality and trend changes, modeled with multiplicative seasonality.
  - **SARIMA**: Incorporates both seasonal and non-seasonal components for autoregressive integrated moving average predictions.
  - **XGBoost**: A tree-based ensemble method optimized for regression tasks, trained with engineered time-series features.

Schedule : Every 1 hour at 50 minutes of each hour, e.g. 01:50, 02:50, 03:50 ... onward. start 2024-12-21 00:00 utc7, end 2024-12-24 06:00 utc7.

*\*More detail please check on code at dag2.py*

## DAG 2 : Train\_Model



```
# Save all models
import pickle
models = {
    'prophet': model,
    'sarima': sarima_results,
    'xgb': xgb_model,
    'features': features,
    'last_timestamp': df['datetime'].max(),
    'last_prices': df['lastPrice'].tail(6).tolist() # Save last 6 prices for lag features
}

modelFile = '/usr/local/airflow/output/models.pkl'
with open(modelFile, 'wb') as f:
    pickle.dump(models, f)
```



## DAG 2 : Train\_Model

### Prep Data

Load past 1 day to current data (288 records)

```
data = ti.xcom_pull(
    task_ids="Get_BTC_from_Mongodb",
    key="data"
)
df = pd.DataFrame(json.loads(data))

# prep data
df['datetime'] = pd.to_datetime(df['closeTime'], unit='ms')
df = df[~df['datetime'].duplicated(keep='first')]
df.set_index('datetime', inplace=True, drop=False)
df = df.sort_index()
df['lastPrice'] = pd.to_numeric(df['lastPrice'], errors='coerce')
df['lastPrice'] = df['lastPrice'].ffill() # Forward fill
```

Assign index from datetime and sort by index.

As requirement :  
X = closeTime only  
Y = lastPrice

### Prophet

datetime = closeTime

```
from prophet import Prophet
# Prepare data for Prophet
prophet_df = pd.DataFrame()
prophet_df['ds'] = df['datetime']
prophet_df['y'] = df['lastPrice']

model = Prophet(
    changepoint_prior_scale=0.001,
    seasonality_prior_scale=0.1,
    seasonality_mode='additive',
    daily_seasonality=True,
    weekly_seasonality=True
)
model.fit(prophet_df)
```

### SARIMA

Use index as datetime

```
from statsmodels.tsa.statespace.sarimax import SARIMAX
# Fit SARIMA model
sarima_model = SARIMAX(
    df['lastPrice'],
    order=(0, 1, 3), # ARIMA parameters (p,d,q)
    seasonal_order=(0, 1, 1, 24), # Seasonal parameters (P,D,Q,s)
)
sarima_results = sarima_model.fit()
```

### XGBoost

#### Feature Engineering for XGBoost

```
# Feature Engineering for Time Series
df['lag1'] = df['lastPrice'].shift(1) # previous 1 price as feature
df['lag2'] = df['lastPrice'].shift(2) # previous 2 price as feature
df['lag3'] = df['lastPrice'].shift(3) # previous 3 price as feature
df['rolling_mean_6'] = df['lastPrice'].rolling(window=6).mean() #
df['rolling_std_6'] = df['lastPrice'].rolling(window=6).std() #
```

```
# Features and Target
features = ['closeTime', 'lag1', 'lag2', 'lag3', 'rolling_mean_6', 'rolling_std_6']
X = df[features].dropna()
y = df['lastPrice'].loc[X.index]

# Train-test split
split_index = int(0.8 * len(X))
X_train = X[:split_index]
X_test = X[split_index:]
y_train = y[:split_index]
y_test = y[split_index:]

print(f"Training set size: {X_train.shape}, Test set size: {X_test.shape}")

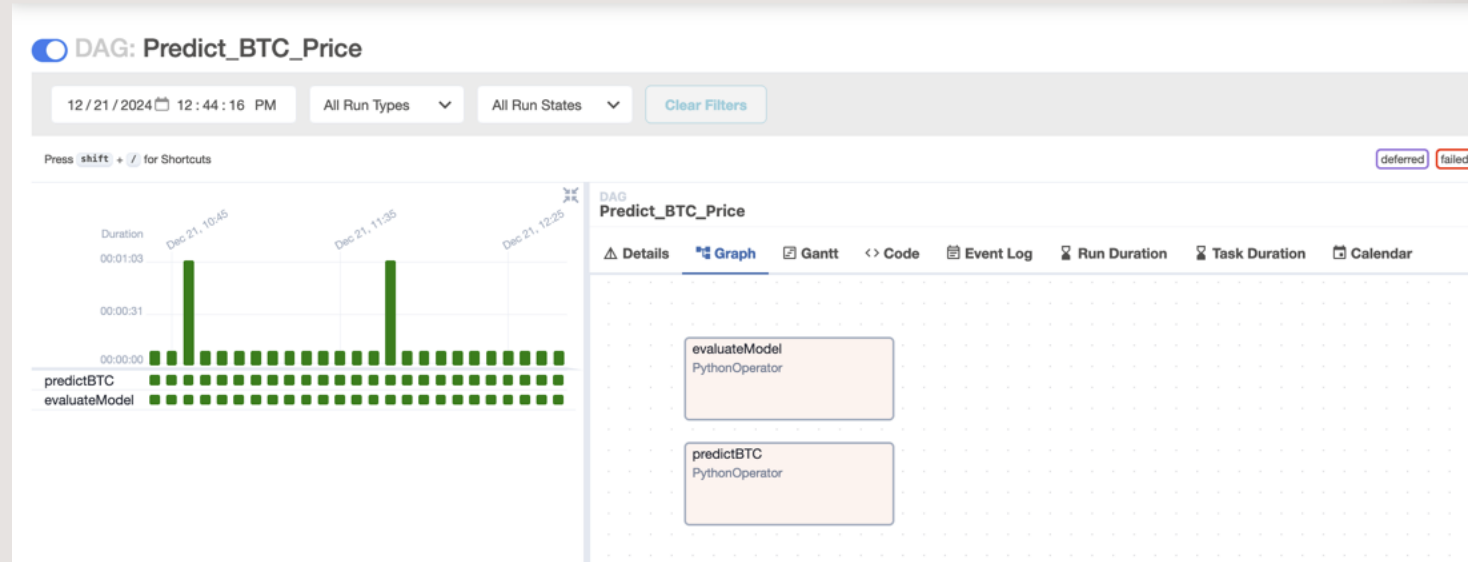
# XGBoost parameters
params = {
    'objective': 'reg:squarederror',
    'learning_rate': 0.01,
    'max_depth': 5,
    'subsample': 0.8,
    'colsample_bytree': 0.8,
    'reg_alpha': 1.0,
    'reg_lambda': 2.0,
    'min_child_weight': 1,
    'eval_metric': 'rmse'
}

# Create DMatrix objects
dtrain = xgb.DMatrix(X_train, label=y_train)
dtest = xgb.DMatrix(X_test, label=y_test)

# Set up evaluation list
evallist = [(dtrain, 'train'), (dtest, 'eval')]

# Train with early stopping
xgb_model = xgb.train(
    params,
    dtrain,
    num_boost_round=1000,
    evals=evallist,
    early_stopping_rounds=50,
    verbose_eval=True
)
```

## DAG 3 : Predict\_BTC\_Price



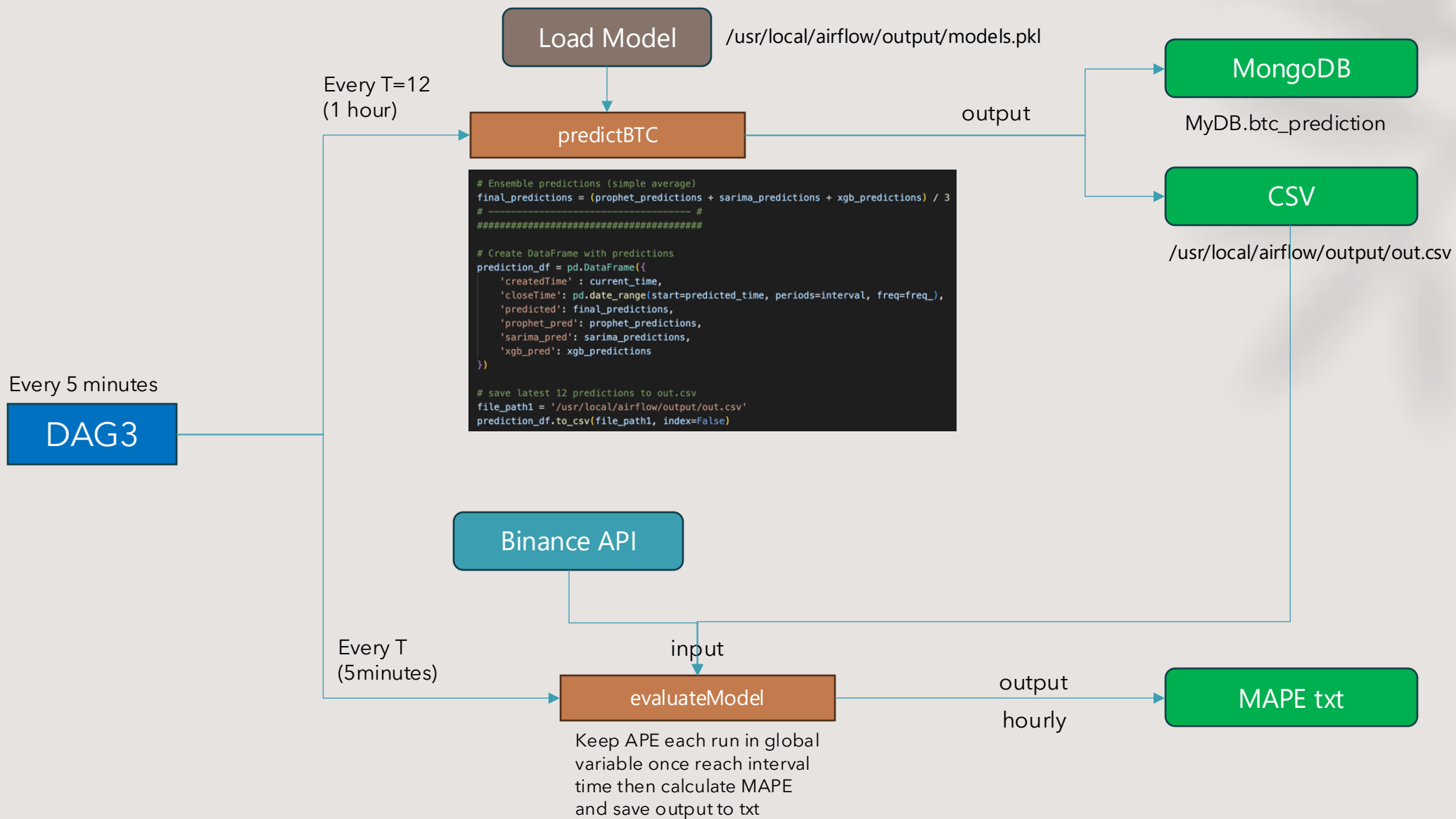
There are 2 independent task in DAG 3 as following :

- **predictBTC** : This task execute only every interval (in this case interval = 12 runs) has been met. So, each run in every 5 minutes means  $12 * 5\text{minutes} = 1\text{hour}$ . It predicts using the saved model from dag2 output and **forecast BTC price in next 12-time interval**, e.g. if task executed at 10:55 then it will forecast BTC price for 11:00, 11:05, 11:10,... 11:55 and then save(overwrite) the latest output to **"/usr/local/airflow/output/out.csv"** as well as append it to the historical to **"/usr/local/airflow/output/prediction\_history.csv"** for later analysis.
- **EvaluateModel** : use data output from above task **"/usr/local/airflow/output/out.csv"** (use ensemble predicted result) to calculate MAPE with actual prices.
  - Every run (5 minutes) it will load current lastPrice from Binance API as y-actual and compare error with y-hat (ensemble predicted result) and add into global variable list.
  - Every interval has been met (12 runs) it will calculate MAPE(Mean Absolute Percentage Error) from last 12 absolute error percentage and save to **f'/usr/local/airflow/output/mape\_{str\_date}.txt'**

Schedule : Every 5 minutes from 22:00-06:00 between december21-24. start 2024-12-21 00:00 utc7, end 2024-12-24 06:00 utc7

*\*More detail please check on code at dag3.py*

## DAG 3 : Predict\_BTC\_Price



# Example of Output

out.csv

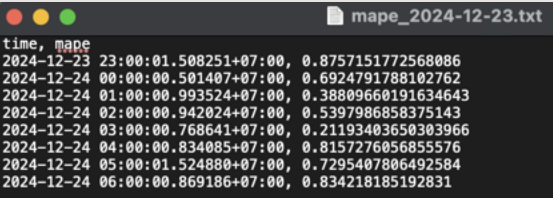
- write every 1 hour

createdTime	closeTime	predicted	prophet_pred	sarima_pred	xgb_pred
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:05:02.328255+07:00	95480.67276	95219.5978	95462.08454	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:10:02.328255+07:00	95507.3187	95236.55105	95525.06911	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:15:02.328255+07:00	95468.04296	95253.84045	95389.95248	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:20:02.328255+07:00	95549.81703	95271.32149	95617.79367	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:25:02.328255+07:00	95497.2787	95288.84884	95442.65131	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:30:02.328255+07:00	95512.43296	95306.27737	95470.68559	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:35:02.328255+07:00	95529.39509	95323.46316	95504.38619	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:40:02.328255+07:00	95495.47118	95340.26451	95385.81309	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:45:02.328255+07:00	95481.42088	95356.5429	95327.3838	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:50:02.328255+07:00	95485.46549	95372.16391	95323.89661	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 06:55:02.328255+07:00	95481.93737	95386.99816	95298.478	95760.336
2024-12-23 06:00:01.328255+07:00	2024-12-23 07:00:02.328255+07:00	95379.66291	95400.92213	94977.73066	95760.336

Use as predicted price each T(5 minutes) to calculate APE and then MAPE at the end of T(12)  
Remark: predicted =(prophet\_pred+sarima\_pred+xgb\_pred)/3

mape\_yyyy-mm-dd.txt

- append every 1 hour

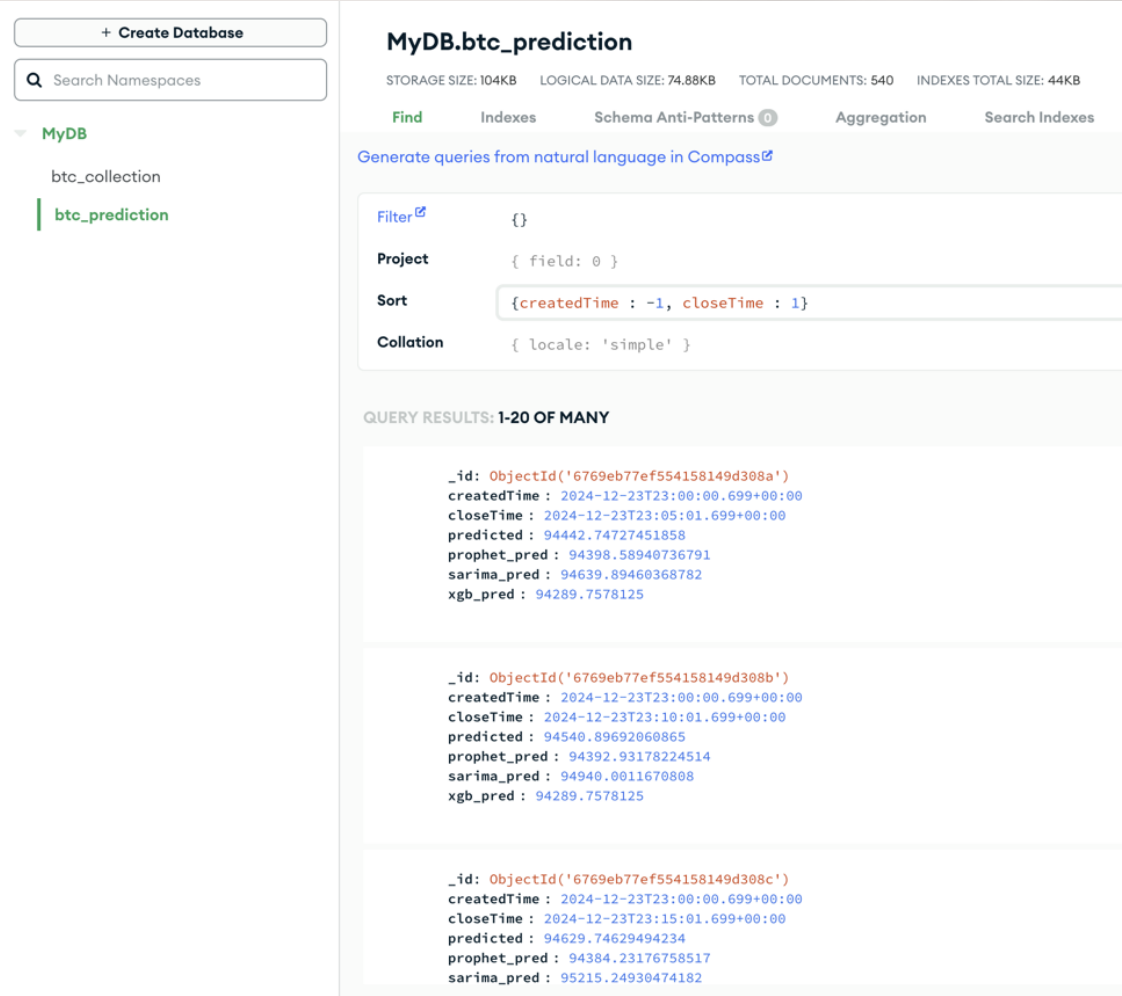


time	mape
2024-12-23 23:00:01.508251+07:00	0.8757151772568086
2024-12-24 00:00:00.501407+07:00	0.6924791788102762
2024-12-24 01:00:00.993524+07:00	0.38809660191634643
2024-12-24 02:00:00.942024+07:00	0.5397986858375143
2024-12-24 03:00:00.768641+07:00	0.21193403650393966
2024-12-24 04:00:00.834085+07:00	0.8157276056855576
2024-12-24 05:00:01.524880+07:00	0.7295407806492584
2024-12-24 06:00:00.869186+07:00	0.834218185192831

Use for calculate MAPE(mean absolute percentage error) from APE of each 5 minutes past hour (12T, 5 minutes each)

## Mongodb btc\_prediction

- Insert the prediction of every 1 hours in advance into mongo for later analysis.



**MyDB.btc\_prediction**

STORAGE SIZE: 104KB LOGICAL DATA SIZE: 74.88KB TOTAL DOCUMENTS: 540 INDEXES TOTAL SIZE: 44KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

Generate queries from natural language in Compass

Filter {}

Project { field: 0 }

Sort {createdTime : -1, closeTime : 1}

Collation { locale: 'simple' }

QUERY RESULTS: 1-20 OF MANY

**Document 1:**

- \_id: ObjectId('6769eb77ef554158149d308a')
- createdTime: 2024-12-23T23:00:00.699+00:00
- closeTime: 2024-12-23T23:05:01.699+00:00
- predicted: 94442.74727451858
- prophet\_pred: 94398.58940736791
- sarima\_pred: 94639.89460368782
- xgb\_pred: 94289.7578125

**Document 2:**

- \_id: ObjectId('6769eb77ef554158149d308b')
- createdTime: 2024-12-23T23:00:00.699+00:00
- closeTime: 2024-12-23T23:10:01.699+00:00
- predicted: 94540.89692060865
- prophet\_pred: 94392.93178224514
- sarima\_pred: 94940.0011670808
- xgb\_pred: 94289.7578125

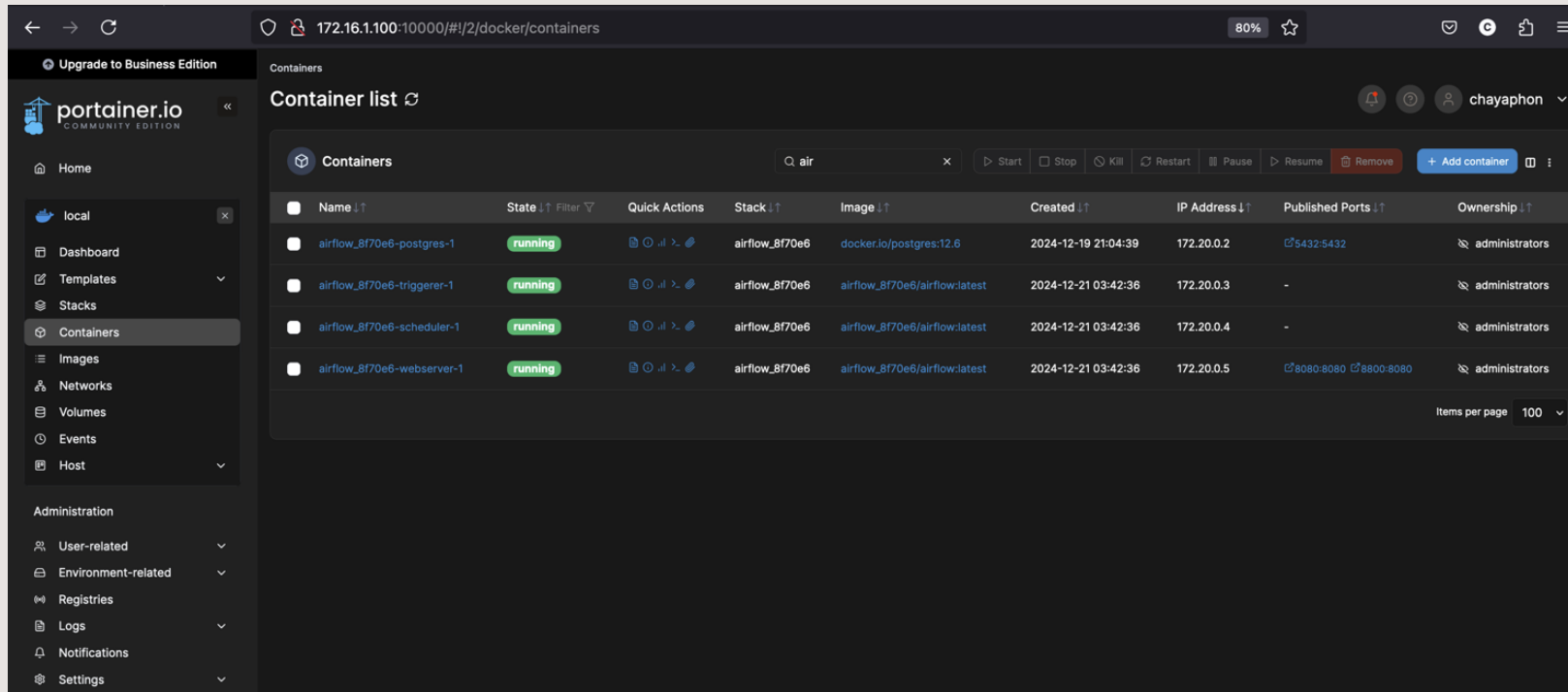
**Document 3:**

- \_id: ObjectId('6769eb77ef554158149d308c')
- createdTime: 2024-12-23T23:00:00.699+00:00
- closeTime: 2024-12-23T23:15:01.699+00:00
- predicted: 94629.74629494234
- prophet\_pred: 94384.23176758517
- sarima\_pred: 95215.24930474182

# Appendix

# Docker

In this project we use docker technology to help us running airflow and related services.



The screenshot displays the Portainer.io web interface, specifically the 'Containers' section. The interface is dark-themed and shows a list of running containers. The left sidebar contains navigation options like Home, local, Dashboard, Templates, Stacks, Containers (selected), Images, Networks, Volumes, Events, Host, and Administration. The main area shows a 'Container list' with a search bar and various action buttons (Start, Stop, Kill, Restart, Pause, Resume, Remove, Add container). The table lists four containers, all in a 'running' state, managed by 'administrators'.

Name	State	Quick Actions	Stack	Image	Created	IP Address	Published Ports	Ownership
airflow_8f70e6-postgres-1	running	[Icons]	airflow_8f70e6	docker.io/postgres:12.6	2024-12-19 21:04:39	172.20.0.2	5432:5432	administrators
airflow_8f70e6-triggerer-1	running	[Icons]	airflow_8f70e6	airflow_8f70e6/airflow:latest	2024-12-21 03:42:36	172.20.0.3	-	administrators
airflow_8f70e6-scheduler-1	running	[Icons]	airflow_8f70e6	airflow_8f70e6/airflow:latest	2024-12-21 03:42:36	172.20.0.4	-	administrators
airflow_8f70e6-webserver-1	running	[Icons]	airflow_8f70e6	airflow_8f70e6/airflow:latest	2024-12-21 03:42:36	172.20.0.5	8080:8080	administrators

## docker-compose.override.yml

Override the standard config of astro docker

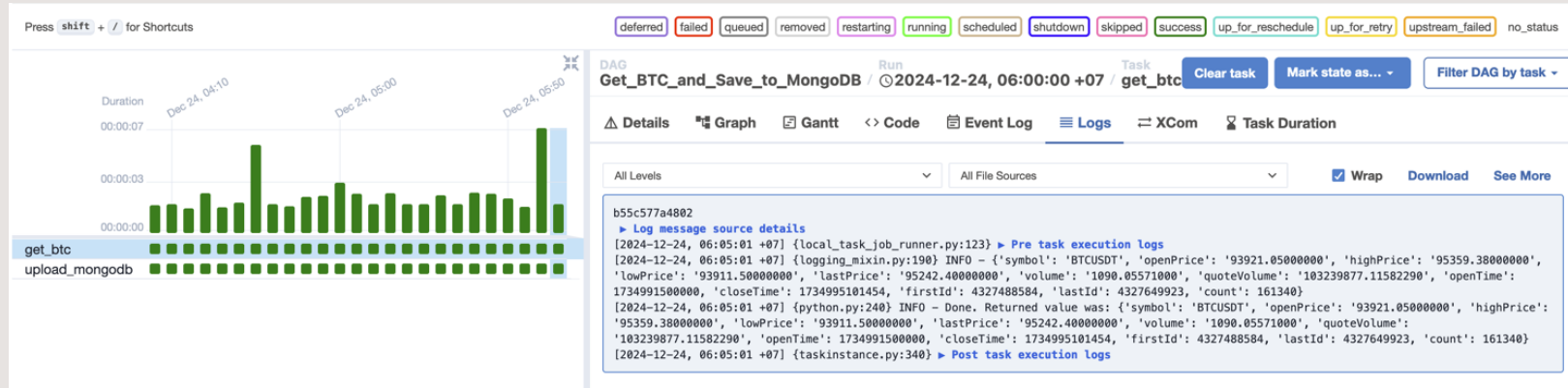
```
version: '3'
services:
  scheduler:
    volumes:
      - ./output:/usr/local/airflow/output
  webserver:
    # Map host port 8080 to container port 8080, Because my docker server has already used port 8080 for other service.
    ports:
      - "8080:8080"
    # below command allow me to access it from client pc (I hosted docker in server and access it from client)
    command: >
      airflow webserver --port 8080 --host 0.0.0.0
```

I have created local directory "**output**" which map to docker airflow "**usr/local/airflow/output**" where I use to dump output of task inside docker to persistent in local machine.

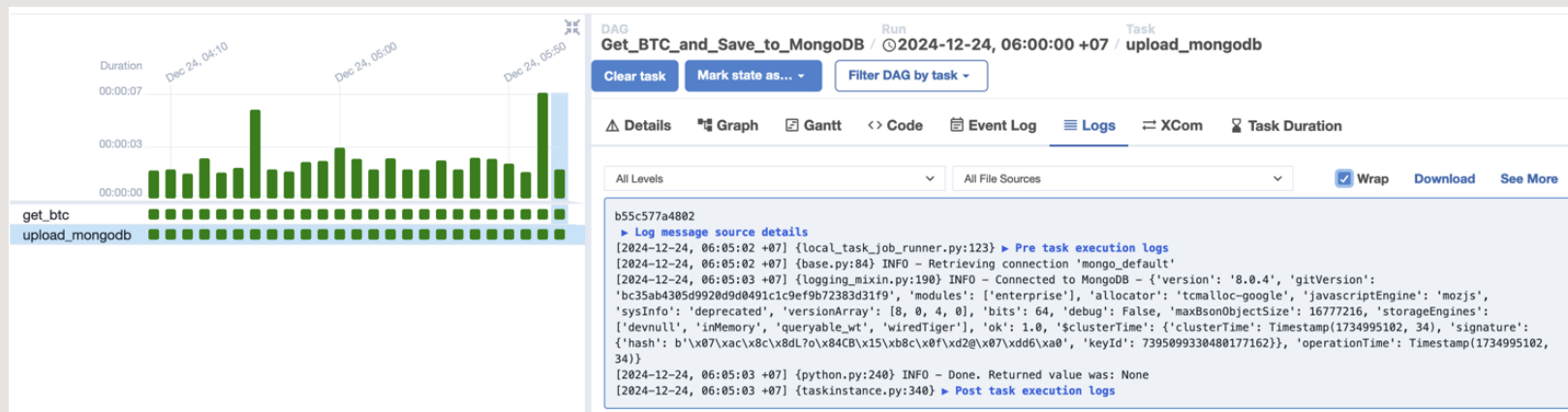
# Some of Task Output Screen Shot

## Get\_BTC\_and\_Save\_to\_MongoDB

get\_btc



upload\_mongodb

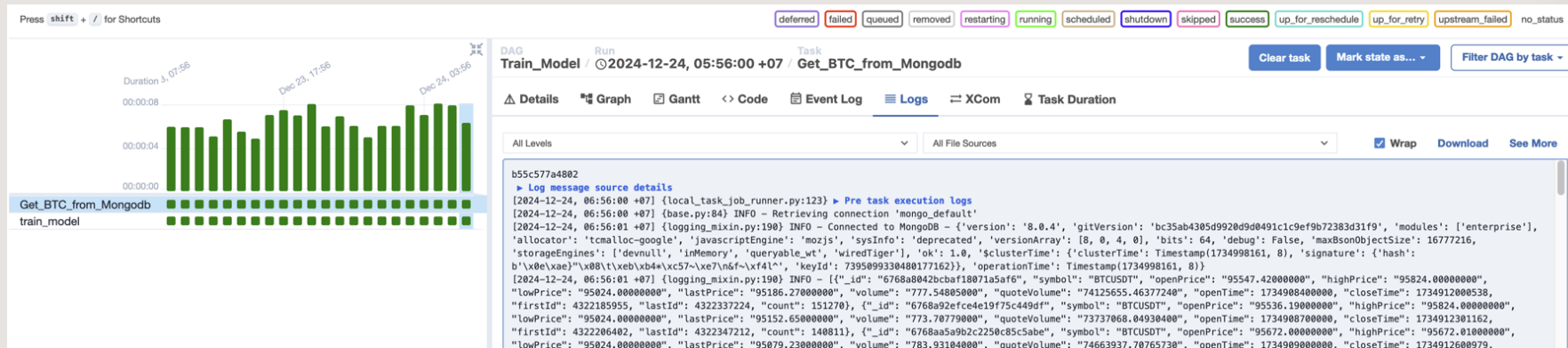




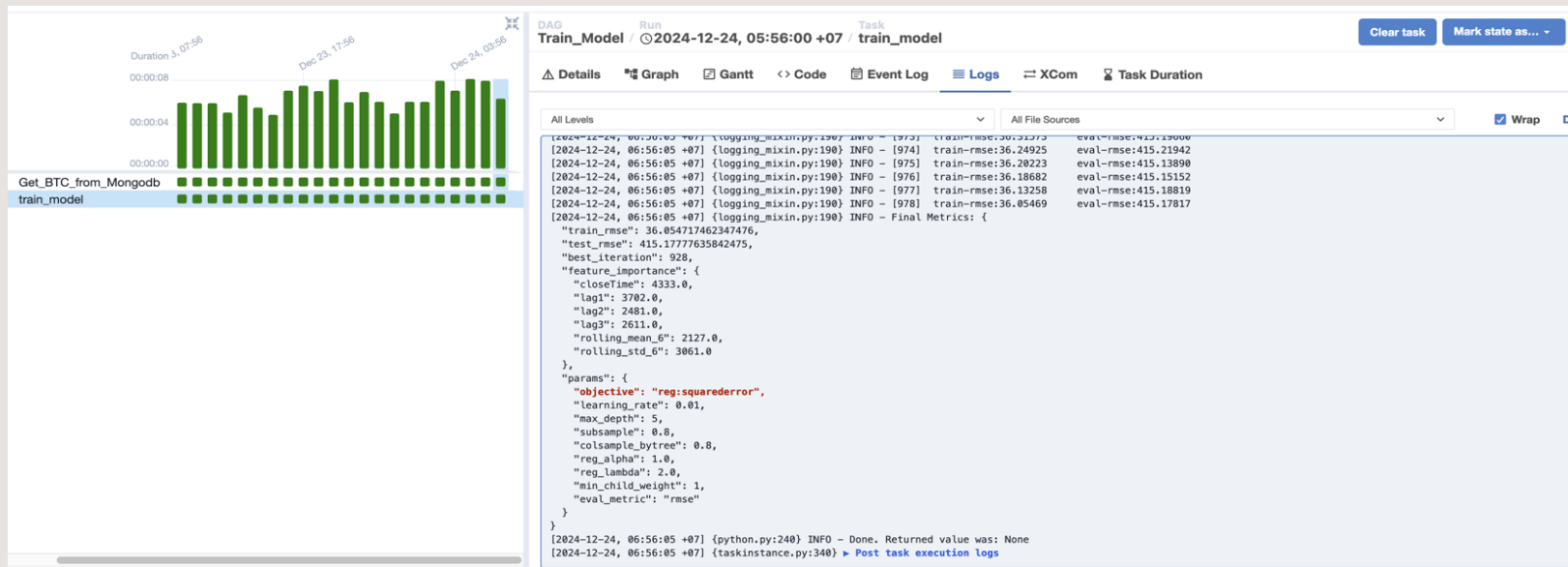
# Some of Task Output Screen Shot

## Tran\_Model

Get\_BTC\_from\_Mongodb. (last 1 day data 288 records)



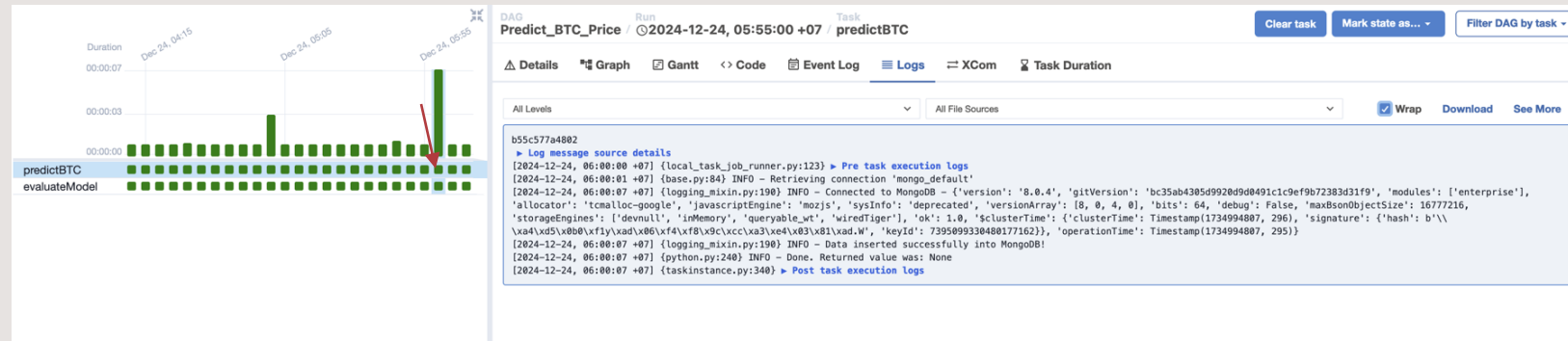
## Train\_model



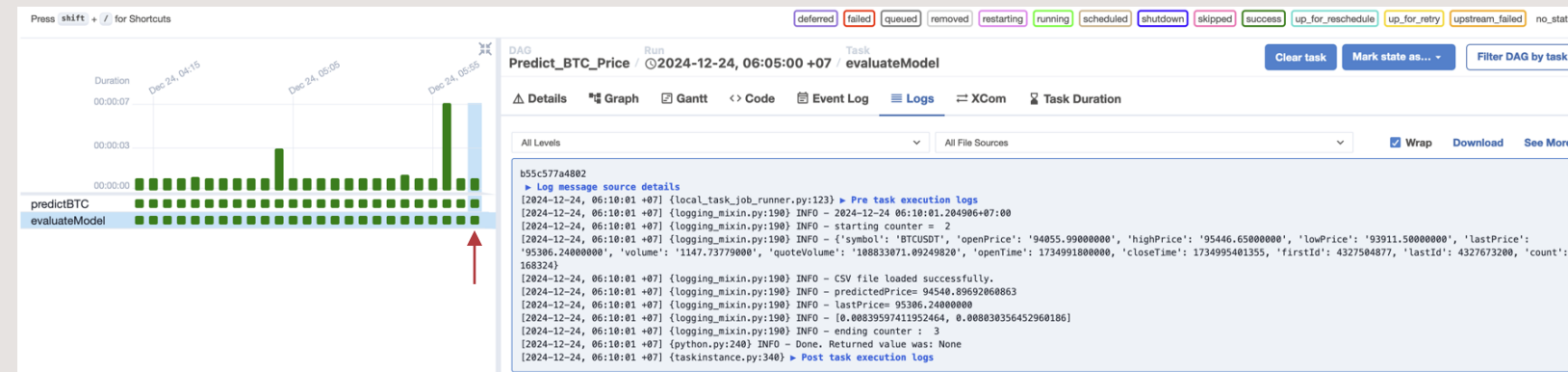
# Some of Task Output Screen Shot

## Tran\_Model

predictBTC : write predicted price in next 12T (5 minutes each) to out.csv and insert data to MongoDB



evaluateModel : calculate error (predict-actual)/actual for each T (5 minutes)



# Mongo DB

btc\_collection

+ Create Database

Q Search Namespaces

MyDB

btc\_collection

btc\_prediction

MyDB.btc\_collection

STORAGE SIZE: 236KB LOGICAL DATA SIZE: 381.83KB TOTAL DOCUMENTS: 1289 INDEXES TOTAL SIZE: 68KB

Find

Indexes

Schema Anti-Patterns 0

Aggregation

Search Indexes

Generate queries from natural language in Compass

Filter {}

Project { field: 0 }

Sort {closeTime : -1}

Collation { locale: 'simple' }

QUERY RESULTS: 1-20 OF MANY

\_id: ObjectId('6769ec9f498a6cb58cc50398')

symbol: "BTCUSD"

openPrice: "93921.05000000"

highPrice: "95359.38000000"

lowPrice: "93911.50000000"

lastPrice: "95242.40000000"

volume: "1090.05571000"

quoteVolume: "103239877.11582290"

openTime: 1734991500000

closeTime: 1734995101454

firstId: 4327488584

lastId: 4327649923

count: 161340

\_id: ObjectId('6769eb767943774dea094bc5')

symbol: "BTCUSD"

openPrice: "93964.01000000"

highPrice: "95191.81000000"

lowPrice: "93911.50000000"

lastPrice: "95191.81000000"

volume: "1019.87653000"

quoteVolume: "96479148.94925870"

openTime: 1734991200000

closeTime: 1734994800641

btc\_prediction

+ Create Database

Q Search Namespaces

MyDB

btc\_collection

btc\_prediction

MyDB.btc\_prediction

STORAGE SIZE: 104KB LOGICAL DATA SIZE: 74.88KB TOTAL DOCUMENTS: 540 INDEXES TOTAL SIZE: 44KB

Find

Indexes

Schema Anti-Patterns 0

Aggregation

Search Indexes

Generate queries from natural language in Compass

Filter {}

Project { field: 0 }

Sort {createdTime : -1, closeTime : 1}

Collation { locale: 'simple' }

QUERY RESULTS: 1-20 OF MANY

\_id: ObjectId('6769eb77ef554158149d308a')

createdTime: 2024-12-23T23:00:00.699+00:00

closeTime: 2024-12-23T23:05:01.699+00:00

predicted: 94442.74727451858

prophet\_pred: 94398.58940736791

sarima\_pred: 94639.89460368782

xgb\_pred: 94289.7578125

\_id: ObjectId('6769eb77ef554158149d308b')

createdTime: 2024-12-23T23:00:00.699+00:00

closeTime: 2024-12-23T23:10:01.699+00:00

predicted: 94540.89692060865

prophet\_pred: 94392.93178224514

sarima\_pred: 94940.0011670808

xgb\_pred: 94289.7578125

\_id: ObjectId('6769eb77ef554158149d308c')

createdTime: 2024-12-23T23:00:00.699+00:00

closeTime: 2024-12-23T23:15:01.699+00:00

predicted: 94629.74629494234

prophet\_pred: 94384.23176758517

sarima\_pred: 95215.24930474182

Dads6005 Final Project Assignment : Chayaphon Sornthananon 6610422007

19

## Code Reference

[https://github.com/chayaphon/realtime\\_airflow-btcPrediction](https://github.com/chayaphon/realtime_airflow-btcPrediction)