

MLOPs Project

Environmental Monitoring and Pollution Prediction System

Muhammad Ali

21I-0887

A

Task 1: Managing Environmental Data with DVC

The task involved setting up an automated DVC pipeline that will automatically get data using IQAir API and then automatically push the data changes to Google Drive which we are using as our remote source.

The automation job is scheduled to run every 10 minute to get data from IQ AIR and then pushes relevant changes to Github and Google Drive.

For Google Drive, a service account was configured that provided us with a JSON file to add to our application for data pushing. All environment variables were kept safely in an environment file.

The Python Script was written that automatically fetches important metrics such as time, city, country, temperature, humidity, wind speed, air quality us, and air quality cn and saves them as CSV and JSON files accordingly. The Python Script is then run by a BAT script every 10 minute to log all data.

The commands used to automate the process are:

```
Windows PowerShell
PS G:\Semester 7\ML0Ps\course-project-ali0887> python fetch_weather_data.py
Data saved successfully:
JSON: data\weather_data_20241214_190058.json
CSV: data\weather_data_20241214_190058.csv
PS G:\Semester 7\ML0Ps\course-project-ali0887> python fetch_weather_data.py
Data saved successfully:
JSON: data\weather_data_20241214_190101.json
CSV: data\weather_data_20241214_190101.csv
PS G:\Semester 7\ML0Ps\course-project-ali0887> python fetch_weather_data.py
Data saved successfully:
JSON: data\weather_data_20241214_190104.json
CSV: data\weather_data_20241214_190104.csv
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc add data/
100% Adding... |1/1 [00:00, 2.15file/s]

To track the changes with git, run:

    git add .gitignore data.dvc

To enable auto staging, run:

    dvc config core.autostage true
PS G:\Semester 7\ML0Ps\course-project-ali0887> git add .gitignore data.dvc
PS G:\Semester 7\ML0Ps\course-project-ali0887> git commit -m "Adding Raw Data"
[main 00622d4] Adding Raw Data
5 files changed, 13 insertions(+)
create mode 100644 .dvc/.gitignore
create mode 100644 .dvc/config
create mode 100644 .dvcignore
create mode 100644 .gitignore
create mode 100644 data.dvc
PS G:\Semester 7\ML0Ps\course-project-ali0887> git push origin main
Enumerating objects: 9, done.
Counting objects: 100% (9/9), done.
Delta compression using up to 8 threads
Compressing objects: 100% (5/5), done.
Writing objects: 100% (8/8), 762 bytes | 762.00 KiB/s, done.
Total 8 (delta 0), reused 0 (delta 0), pack-reused 0
To https://github.com/NUCES-ISB/course-project-ali0887.git
be1e78b..00622d4 main -> main
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc remote add -d myremote gdrive://10CtnFhHMNmmlr1XkFGjF_rUMc0Y-2eZ29

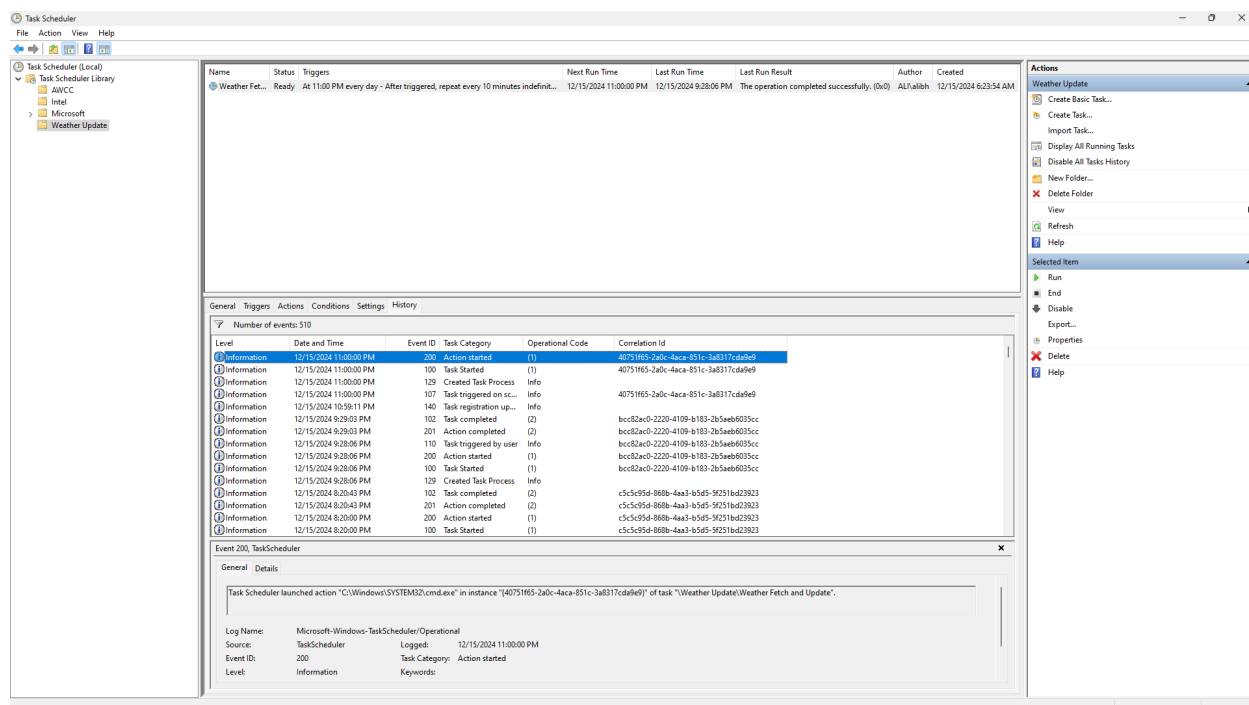
Windows PowerShell
PS G:\Semester 7\ML0Ps\course-project-ali0887> python fetch_weather_data.py
Data saved successfully:
JSON: data\weather_data_20241214_190104.json
CSV: data\weather_data_20241214_190104.csv
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc add data/
100% Adding... |1/1 [00:00, 2.15file/s]

To track the changes with git, run:

    git add .gitignore data.dvc

To enable auto staging, run:

    dvc config core.autostage true
PS G:\Semester 7\ML0Ps\course-project-ali0887> git add .gitignore data.dvc
PS G:\Semester 7\ML0Ps\course-project-ali0887> git commit -m "Adding Raw Data"
[main 00622d4] Adding Raw Data
5 files changed, 13 insertions(+)
create mode 100644 .dvc/.gitignore
create mode 100644 .dvc/config
create mode 100644 .dvcignore
create mode 100644 .gitignore
create mode 100644 data.dvc
PS G:\Semester 7\ML0Ps\course-project-ali0887> git push origin main
Enumerating objects: 9, done.
Counting objects: 100% (9/9), done.
Delta compression using up to 8 threads
Compressing objects: 100% (5/5), done.
Writing objects: 100% (8/8), 762 bytes | 762.00 KiB/s, done.
Total 8 (delta 0), reused 0 (delta 0), pack-reused 0
To https://github.com/NUCES-ISB/course-project-ali0887.git
be1e78b..00622d4 main -> main
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc remote add -d myremote gdrive://10CtnFhHMNmmlr1XkFGjF_rUMc0Y-2eZ29
Setting 'myremote' as a default remote.
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc remote modify myremote gdrive_use_service_account true
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc remote modify myremote gdrive_service_account_json_file_path mlops-final-project-444
714-c95878d44542f.json
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc remote modify myremote gdrive_acknowledge_abuse true
PS G:\Semester 7\ML0Ps\course-project-ali0887> python -m dvc push
10 files pushed
PS G:\Semester 7\ML0Ps\course-project-ali0887> |
```



Task 2: Pollution Trend Prediction with MLFlow

For the task of incorporating MLFlow with our DVC pipeline and then filtering out the best model, an MLFlow experiment was setup that would test and log a combination of models with the following parameters:

Sequence Length : 2, 6, 12, 24

Hidden LSTM Units: 16, 32, 64, 128

Dropout: 0.1, 0.2, 0.3

Learning Rate: 0.001

Epochs: 30, 50

Batch Size: 4, 8, 16

The training is done in Tensorflow using 2 LSTM layers with a dropout layer after each and a Dense layer for predicting.

The total number of combinations came out to 192 and all the results were logged in MLFlow. Chosen Metrics are as follow:

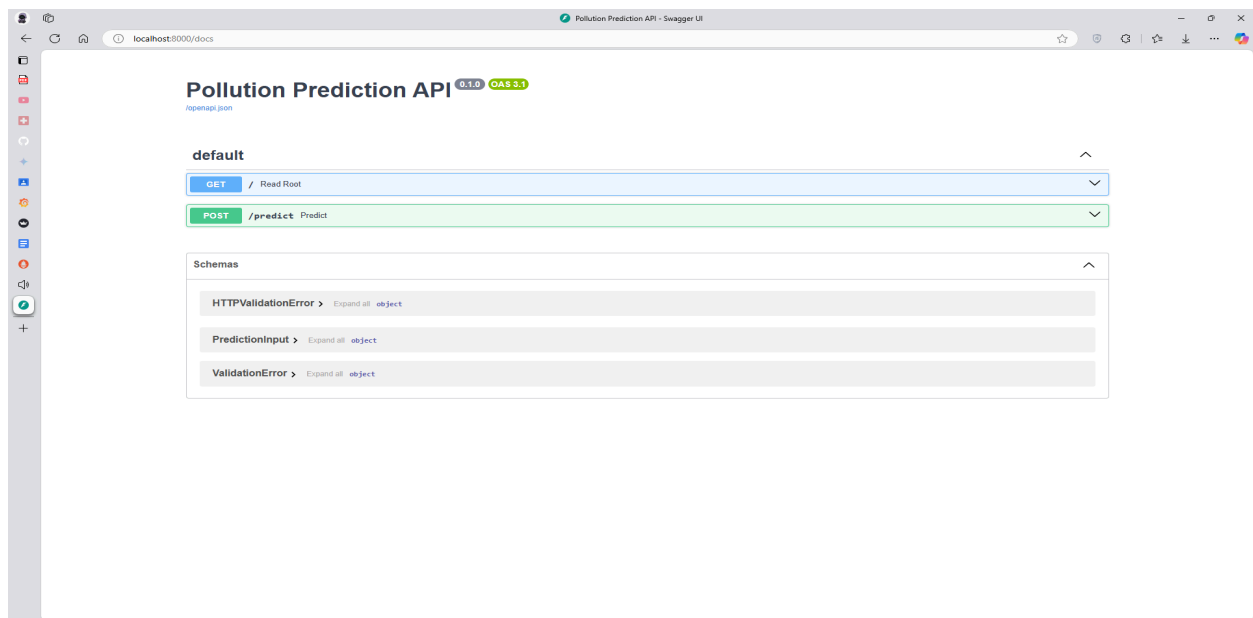
1. Mean Accuracy Error (MAE)
2. Mean Squared Error (MSE)
3. Root Mean Squared Error (RMSE)

As per the results, the best results came for the model with sequence length 2, with 128 hidden lstm units which makes sense since in a timeframe of 20 minutes results will not change much and 128 hidden lstm units allow for a better inference due to more lstm cell retaining more knowledge.

However, during inference the best results came with the model with sequence length = 24 or 4 hours history.

This could be partially because more data equals better predictions. Once again, with 128 hidden lstm units the metrics were the best.

Each Model was deployed using a FAST API on the local machine.



```
Best parameters and RMSE for each sequence length:

Sequence Length 2:
Best parameters: {'sequence_length': 2, 'lstm_units': 128, 'dropout': 0.2, 'learning_rate': 0.001, 'epochs': 50, 'batch_size': 16}
Best RMSE: 0.1758

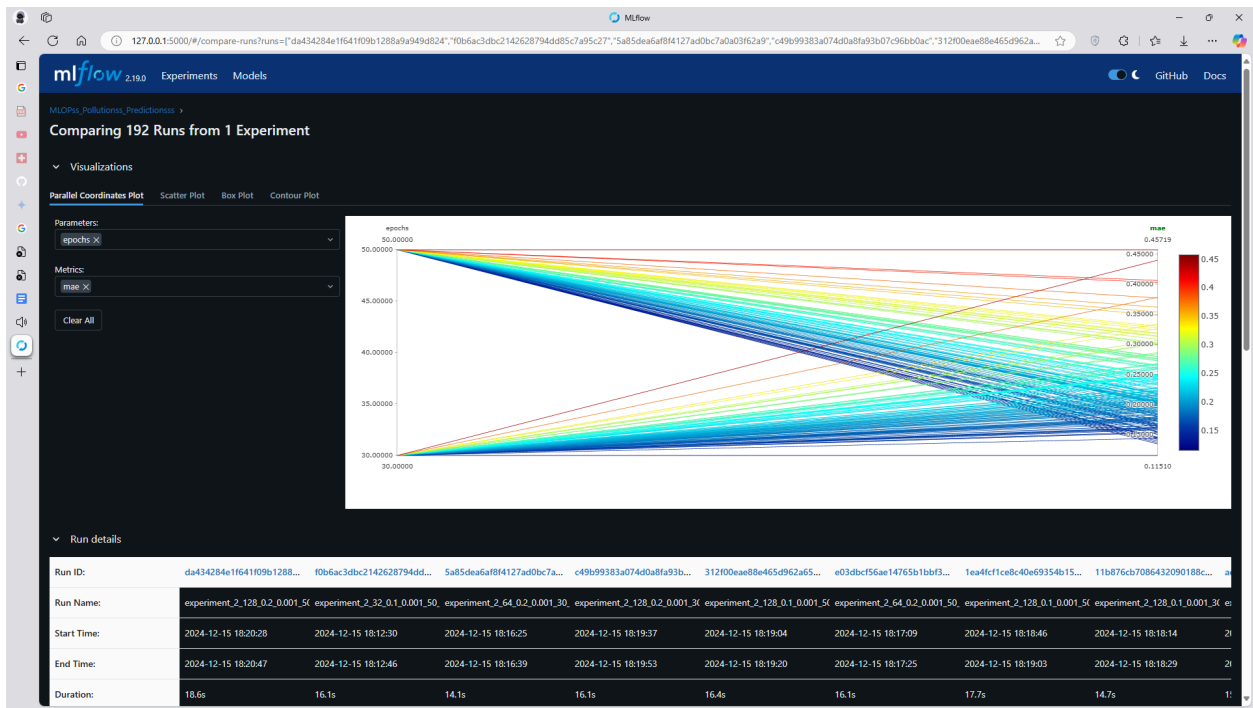
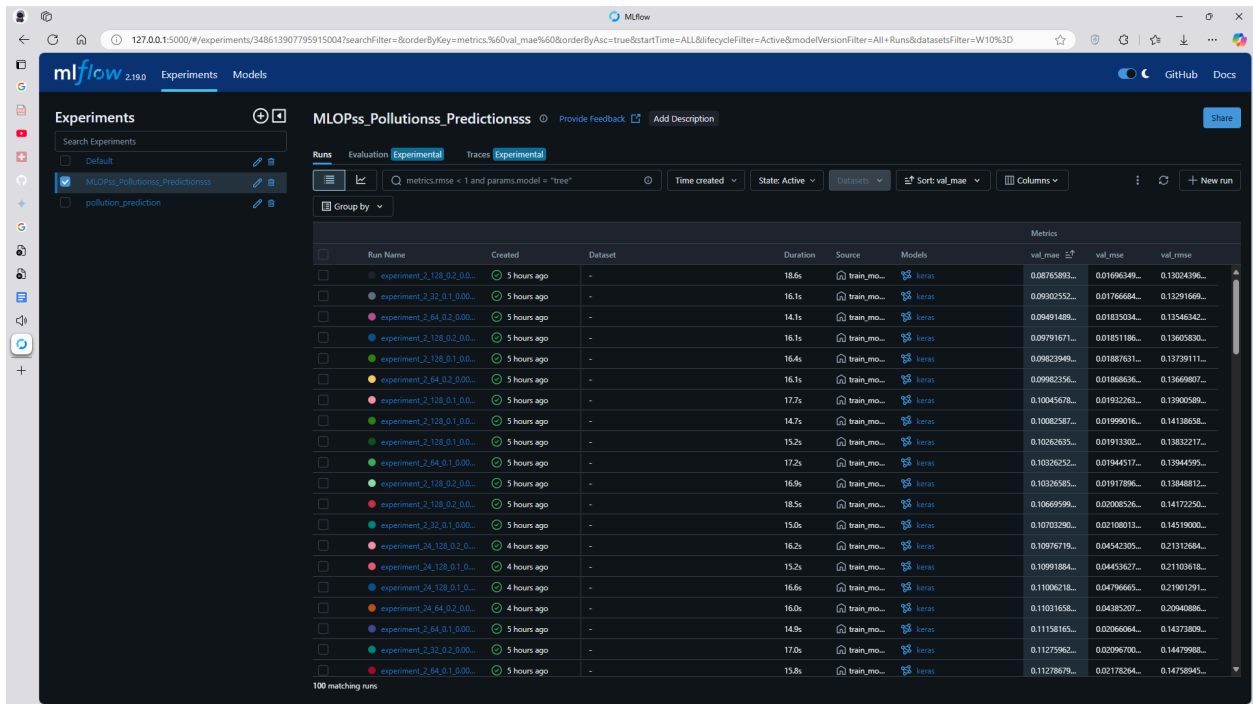
Sequence Length 6:
Best parameters: {'sequence_length': 6, 'lstm_units': 16, 'dropout': 0.2, 'learning_rate': 0.001, 'epochs': 30, 'batch_size': 16}
Best RMSE: 0.2137

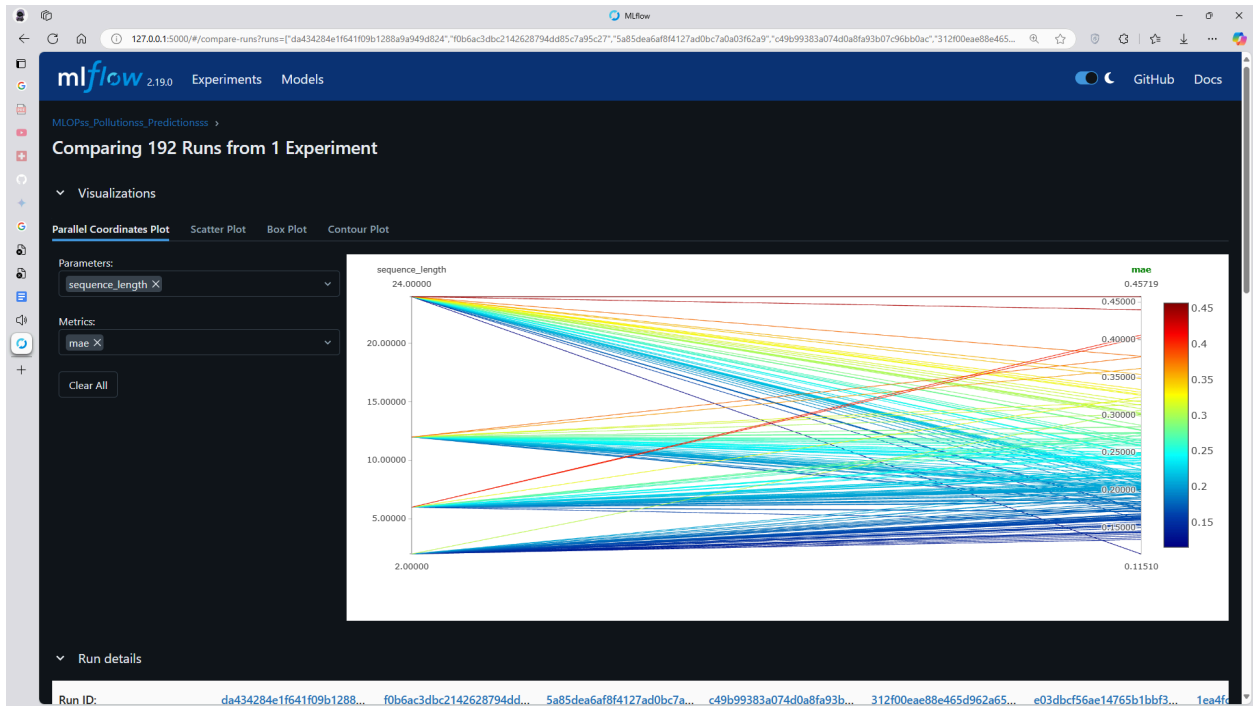
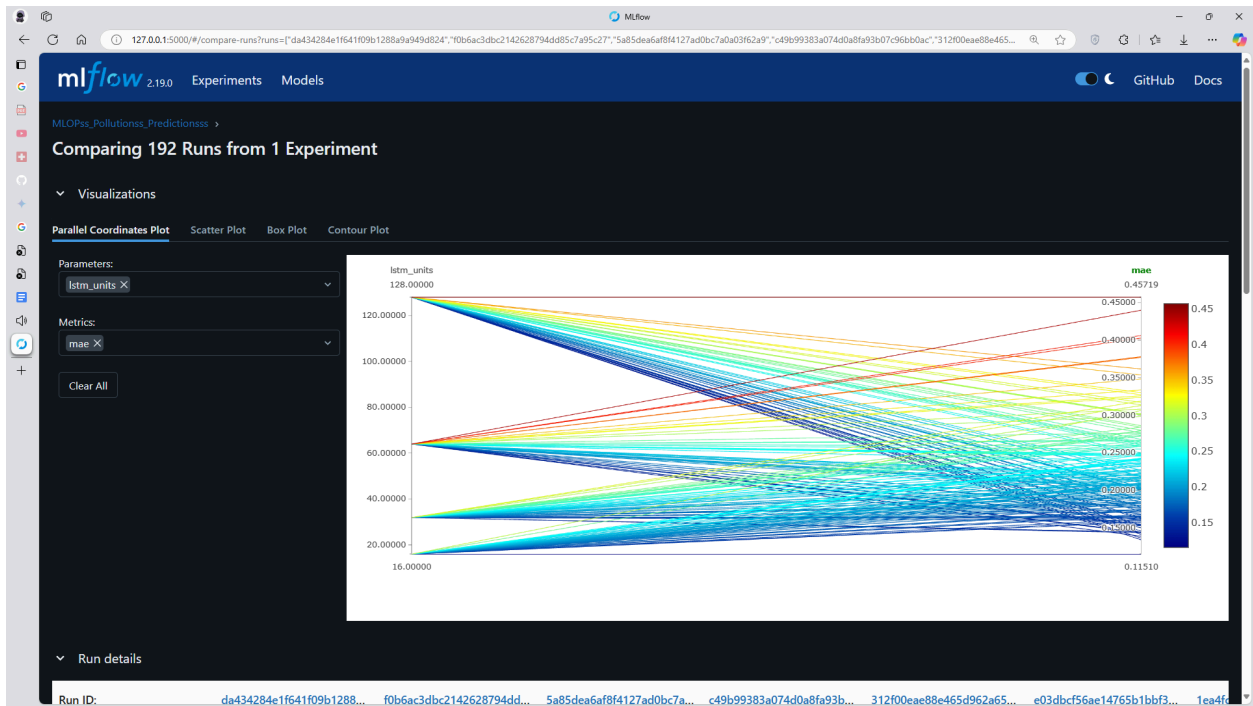
Sequence Length 12:
Best parameters: {'sequence_length': 12, 'lstm_units': 128, 'dropout': 0.2, 'learning_rate': 0.001, 'epochs': 30, 'batch_size': 16}
Best RMSE: 0.2402

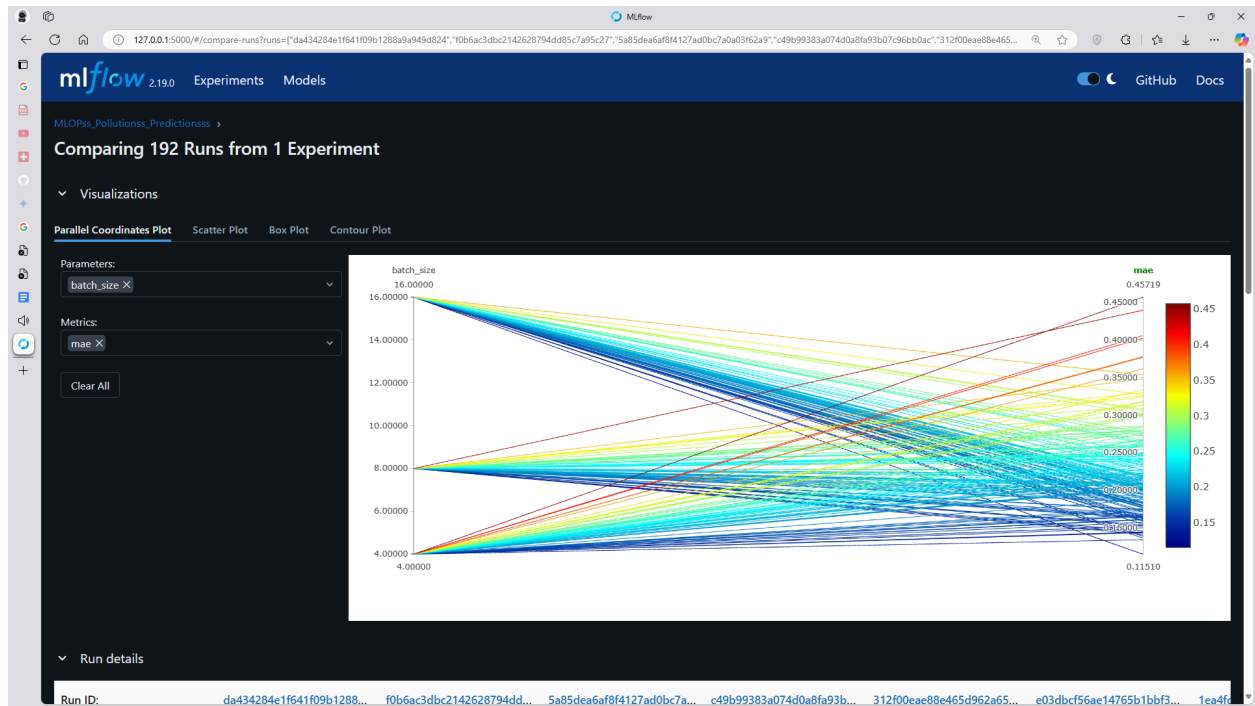
Sequence Length 24:
Best parameters: {'sequence_length': 24, 'lstm_units': 16, 'dropout': 0.1, 'learning_rate': 0.001, 'epochs': 30, 'batch_size': 16}
Best RMSE: 0.1242

Overall best sequence length: 24
Overall best RMSE: 0.1242

Final models have been saved in the 'models' directory:
- models/best_model_seq_2.keras
- models/best_model_seq_6.keras
- models/best_model_seq_12.keras
- models/best_model_seq_24.keras
PS G:\Semester 7\ML\MLPs\course-project-ali0887>
```

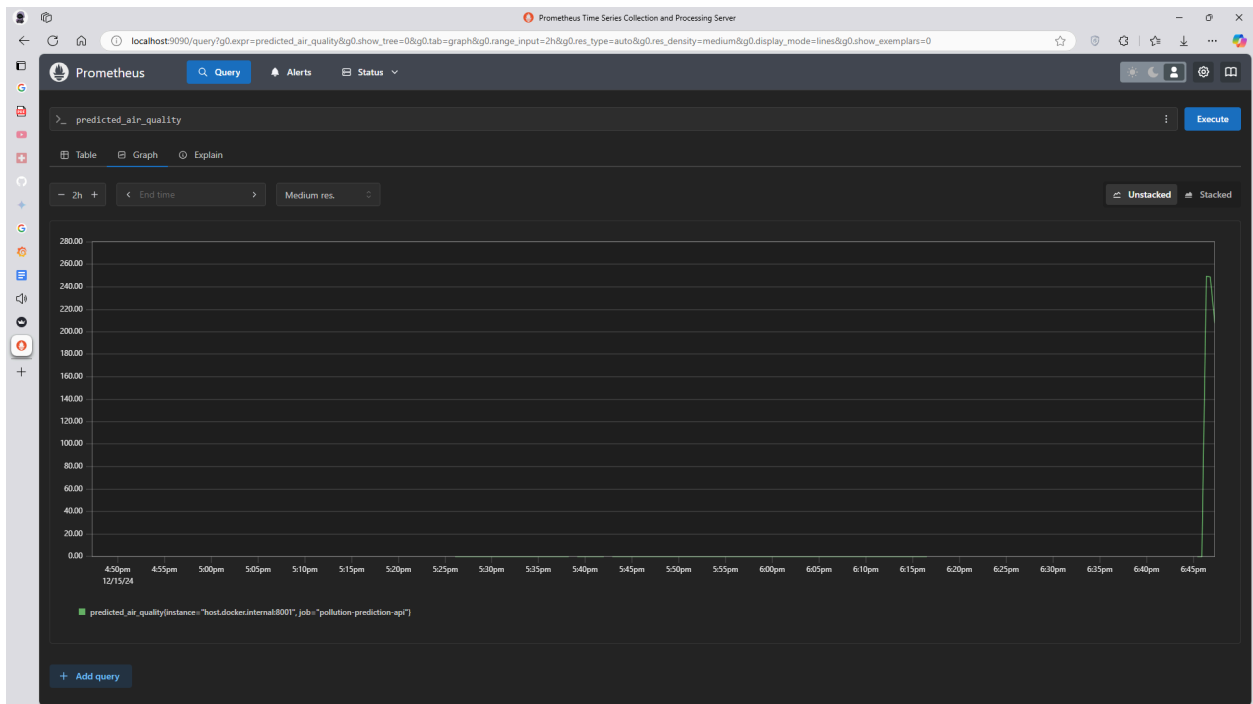
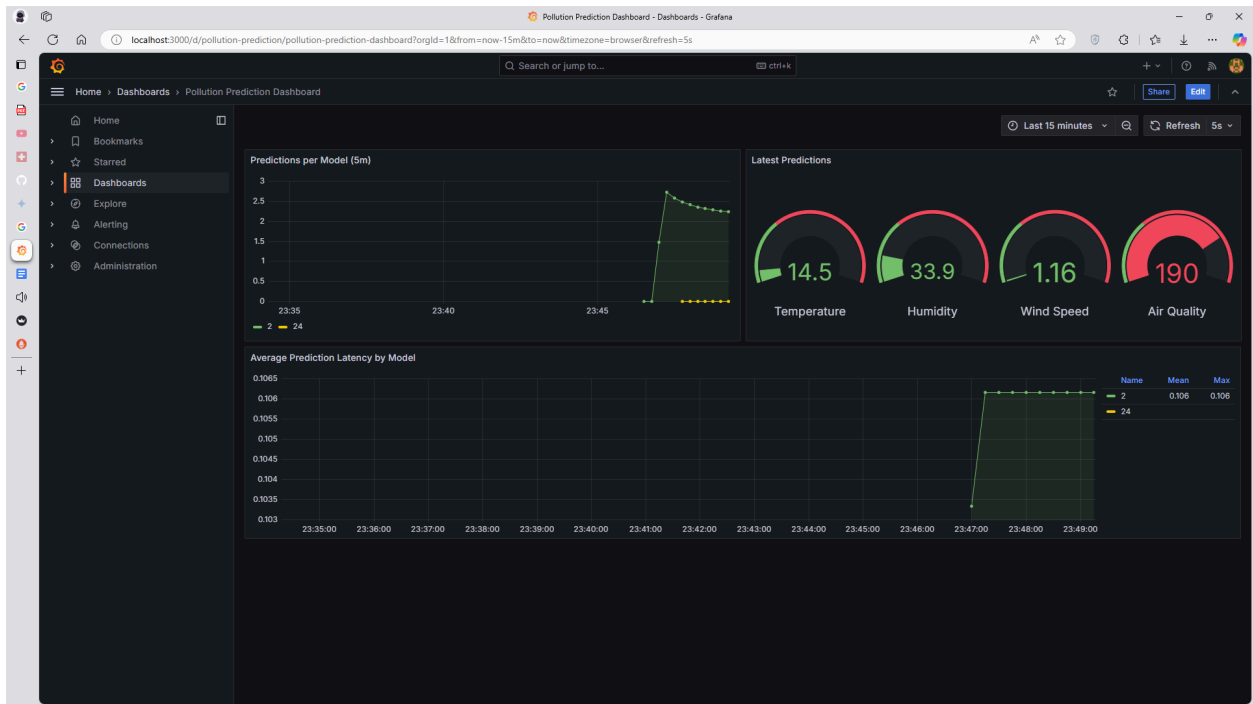






Task 3: Monitoring and Live Testing

Setting up Grafana and Prometheus on Docker involved writing up a Docker Compose Script that automatically fetches relevant data from when the API is called and displays these graphs in Prometheus and Grafana.



Improvements / Optimizations:

As far as Improvements and Optimizations are considered, while not implemented, the identified problems are there not being a model trained on time of day as during the night the temperature usually decreases alongside the overall pollution.

Additionally, incorporating situational awareness into the model like what kind of weather helps to predict the pollution levels better e.g. when it is raining, air pollution is very low.

Moreover, the data collection for this task was low so the model may not converge effectively for different scenarios. Training on more data will help the model identify pollution better.

Alongside there, one may also consider the model architecture as 2 LSTM Layers may not be enough to model a complex task such as weather pollution prediction.