# INSTITUTE FOR ADVANCED COMPUTING
# AND SOFTWARE DEVELOPMENT

AKURDI, PUNE – 411044

Documentation On

**"Metro Failure Detection"**

PG-DBDA FEB 2025

**Submitted By:**

**Group No: 19**

**Pravin Bodke (252535)**
**Suraj Suryavanshi (252549)**

**Mrs. Priti Take**           **Mr. Prashant deshpade**

**Project Guide**           **Centre Coordinator**

# ACKNOWLEDGEMENT

I would like to express my sincere gratitude to Mrs. Priti Take, Project Guide, for providing me with the guidance and support to complete this academic project. her valuable insights , expertise and encouragement have been instrumental in the success of this project.

I would also like to thank my fellow classmates for her support and cooperation during the project. Their feedback and suggestions were helpful in improving the quality of the project.

I would like to extend my gratitude to Mr. Prashant Deshpande , Centre Coordinator, for providing me with the necessary   resources and facilities to complete this project . her support has been crucial in the timely completion of this project.

Finally, I would like to thank my family and friends for her constant encouragement and support throughout the project. her belief in me has been a constant source of motivation and inspiration.

Thank you all for your support and guidance in completing this academic project.

# ABSTRACT

The utilization of predictive maintenance techniques aids in the early identification of machinery failures before they escalate to critical levels. This research introduces a data-driven predictive maintenance framework designed for the air production unit (APU) system within a Metro do Porto train. The framework employs deep learning, specifically autoencoder networks, to efficiently discern abnormal data, thereby significantly reducing false alarms. Numerous analog and digital sensors embedded in the APU system facilitate the detection of behavioral changes and deviations from the norm through data analysis.

Two versions of autoencoders — sparse and variational were implemented, each with both analog and digital sensor data (separately). Experimental findings reveal that digital sensor data is significantly more effective in detecting anomalies, particularly with sparse autoencoders. Performance assessments of the SAE network, using digital sensor data, demonstrate marked improvements in F1 Score, Recall, and Precision — each of which touch almost 100%, when anomalies are flagged in a 24-hour period prior to the failure.

# DECLARATION

I, the undersigned, hereby declare that the project report titled "Predictive Maintenance of Metro Train Air Compressors using Machine Learning and Big Data Technologies" submitted to the Institute For Advanced Computing And Software Development, Akurdi Pune, in partial fulfillment of the requirements for the Post Graduate Diploma in Big Data Analytics (PG DBDA), is my original work.

The project was developed using Python, Apache Kafka, PySpark, and machine learning libraries such as TensorFlow/Keras for modeling. All data preprocessing, model training, and evaluation were carried out by me. The work does not contain any plagiarized content, and no part of it has been submitted earlier for any academic or professional degree. Proper attribution has been given wherever necessary.

Place:                                    Signature:

Date:                                     Name: Pravin Bodke / Suraj Suryavanshi

# Contents

# 1.Introduction

## 1.1 Problem Statement

Unexpected failures of the Air Production Unit (APU) in metro trains disrupt services and increase costs. A predictive maintenance system using real-time sensor data is essential to detect faults early, prevent breakdowns, and enhance reliability while reducing expenses.

## 1.2 Product Scope

The Metro Fault Detection System is developed to monitor and predict equipment failures in metro train air compressors using real-time sensor data streams. It ingests continuous data via Apache Kafka, performs scalable preprocessing using PySpark, and applies a pre-trained LSTM deep learning model to forecast potential faults or maintenance needs.

The system classifies equipment status into "Maintenance Required" or "No Maintenance", enabling timely alerts that can be integrated with operational dashboards or automated notification services. Designed for both batch and streaming modes, the solution supports high-throughput environments and large-scale datasets, making it fully compatible with modern big data ecosystems. This end-to-end pipeline leverages cloud and on-premise technologies to ensure robustness, scalability, and real-time decision-making capabilities for predictive maintenance in metro systems.

### 1.3 Aims & Objectives

The primary goal of this project is to identify and forecast potential faults in metro systems using real-time sensor data. The process begins with collecting and preprocessing large volumes of operational data from the metro network through Kafka and PySpark pipelines. Various machine learning models, with a focus on LSTM for time-series analysis, will be trained on historical data to recognize patterns that indicate normal or faulty operations. After training, the models will be evaluated using accuracy, precision, recall, and F1-score, and the best-performing model will be selected. This chosen model will then be deployed in a real-time environment to predict upcoming failures, enabling proactive maintenance scheduling.

## 2. Overall Description

### 1.4 Workflow of Project:

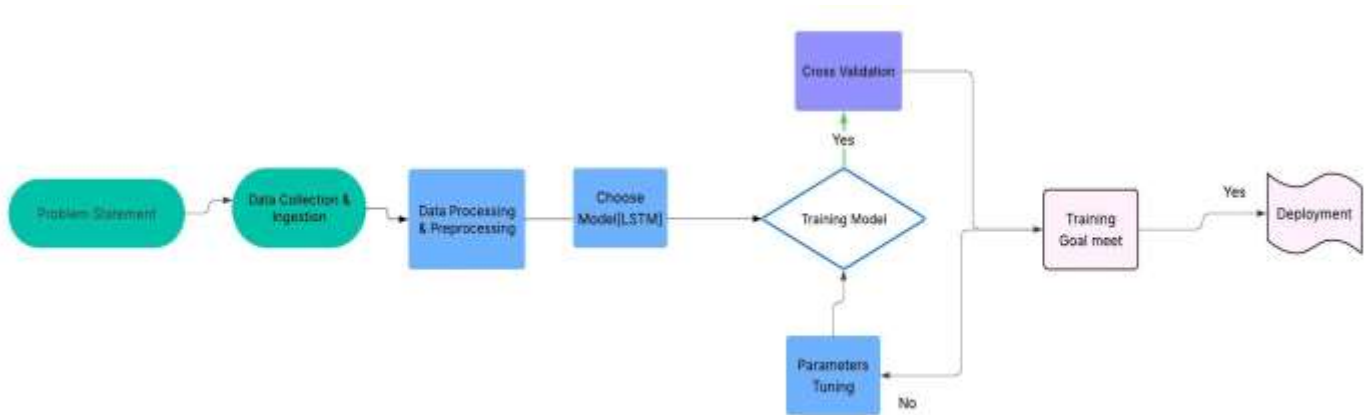The diagram below shows the workflow of this project.



*Figure 1 : Workflow Diagram*

### *Analog Signals:*

1.  *TP2 (bar) — the measure of the pressure on the compressor.*

2.  *TP3 (bar) — the measure of the pressure generated at the pneumatic panel.*

3.  *H1 (bar) — the measure of the pressure generated due to pressure drop when the discharge of the cyclonic separator filter occurs.*

4.  *DV pressure (bar) — the measure of the pressure drop generated when the towers discharge air dryers; a zero reading indicates that the compressor is operating under load.*

5. *Reservoirs (bar) — the measure of the downstream pressure of the reservoirs, which should be close to the pneumatic panel pressure (TP3).*

6. *Motor Current (A) — the measure of the current of one phase of the three-phase motor; it presents values close to 0A — when it turns off, 4A — when working offloaded, 7A — when working under load, and 9A — when it starts working.*

7. *Oil Temperature (°C) — the measure of the oil temperature on the compressor.*

## *Digital Signals:*

1. *COMP — the electrical signal of the air intake valve on the compressor; it is active when there is no air intake, indicating that the compressor is either turned off or operating in an offloaded state.*

2. *DV electric — the electrical signal that controls the compressor outlet valve; it is active when the compressor is functioning under load and inactive when the compressor is either off or operating in an offloaded state.*

3. *TOWERS — the electrical signal that defines the tower responsible for drying the air and the tower responsible for draining the humidity removed from the air; when not active, it indicates that tower one is functioning; when active, it indicates that tower two is in operation.*

4. *MPG — the electrical signal responsible for starting the compressor under load by activating the intake valve when the pressure in the air production unit (APU) falls below 8.2 bar; it activates the COMP sensor, which assumes the same behaviour as the MPG sensor.*

5. *LPS — the electrical signal that detects and activates when the pressure drops below bar*

6. *Pressure Switch — the electrical signal that detects the discharge in the air-drying towers.*

7. *Oil Level — the electrical signal that detects the oil level on the compressor; it*

*is active whe n the oil is below the expected values.*

8.  *Caudal Impulse — the electrical signal that counts the pulse outputs generated by the absolute amount of air flowing from the APU to the reservoirs.*

## 1.5 Data Preprocessing and Cleaning:

**Data Cleaning:**

The raw metro operational data collected from sensors and logs may contain noise, missing values, duplicate records, and irrelevant features. To ensure accurate and reliable predictions, dat aning is carried out as follows:a cle

1.  Handling Missing Values

    Sensor readings may have gaps due to transmission errors or sensor malfunctions. Missing values are handled using statistical imputation methods (mean/median filling) or by forward/backward filling in the time series, ensuring the continuity of the data.

2.  Removing Irrelevant Features

    Some columns in the dataset may not contribute to fault detection (e.g., redundant IDs or static values). Such attributes are dropped to improve model efficiency and reduce complexity.

3.  Outlier Detection and Removal

    Extreme abnormal readings caused by faulty sensors or logging errors are identified using statistical thresholds (e.g., Z-score, IQR method) and are either removed or corrected.

4.  Timestamp Formatting

    All time-related attributes are converted into a standardized datetime format to enable proper time-series sequencing and feature extraction.

**Label encoding:**

In the metro fault detection dataset, the target variable status indicates the operational condition of the system, where the values are in a human-readable format (e.g., *Normal*, *Faulty*). Since machine learning algorithms require numerical input, these categorical labels are converted into numeric form using Label Encoding.

## 1.6 Exploratory Data Analysis:

Exploratory Data Analysis (EDA) is the process of examining the metro fault dataset to identify patterns, trends, and anomalies. Time-series plots were used to observe sensor behavior over different months and detect seasonal variations in faults. Status distribution analysis helped identify class imbalance between normal and fault cases. Correlation analysis revealed relationships between sensor readings and fault occurrence. These insights guided feature selection and model preparation for accurate fault prediction

**Metro Train Failure and Maintenance Schedule** :-

| Nr. | Start Time | End Time | Failure | Severity | Report |
|-----|-----------|----------|---------|----------|--------|
| #1 | 4/18/2020 0:00 | 4/18/2020 23:59 | **Air leak** | High stress | |
| #1 | 5/29/2020 23:30 | 5/30/2020 6:00 | **Air Leak** | High stress | Maintenance on 30Apr at 12:00 |
| #3 | 6/5/2020 10:00 | 6/7/2020 14:30 | **Air Leak** | High stress | Maintenance on 8Jun at 16:00 |
| #4 | 7/15/2020 14:30 | 7/15/2020 19:00 | **Air Leak** | High stress | Maintenance on 16Jul at 00:00 |

*Figure 2: Historical Records of Air Leak Failures with Corresponding Maintenance Actions*

- *This table lists recorded air leak incidents in 2020, along with their start and end times, severity, and scheduled maintenance reports. All events were categorized as "High stress" failures, requiring prompt maintenance interventions to restore system integrity.*

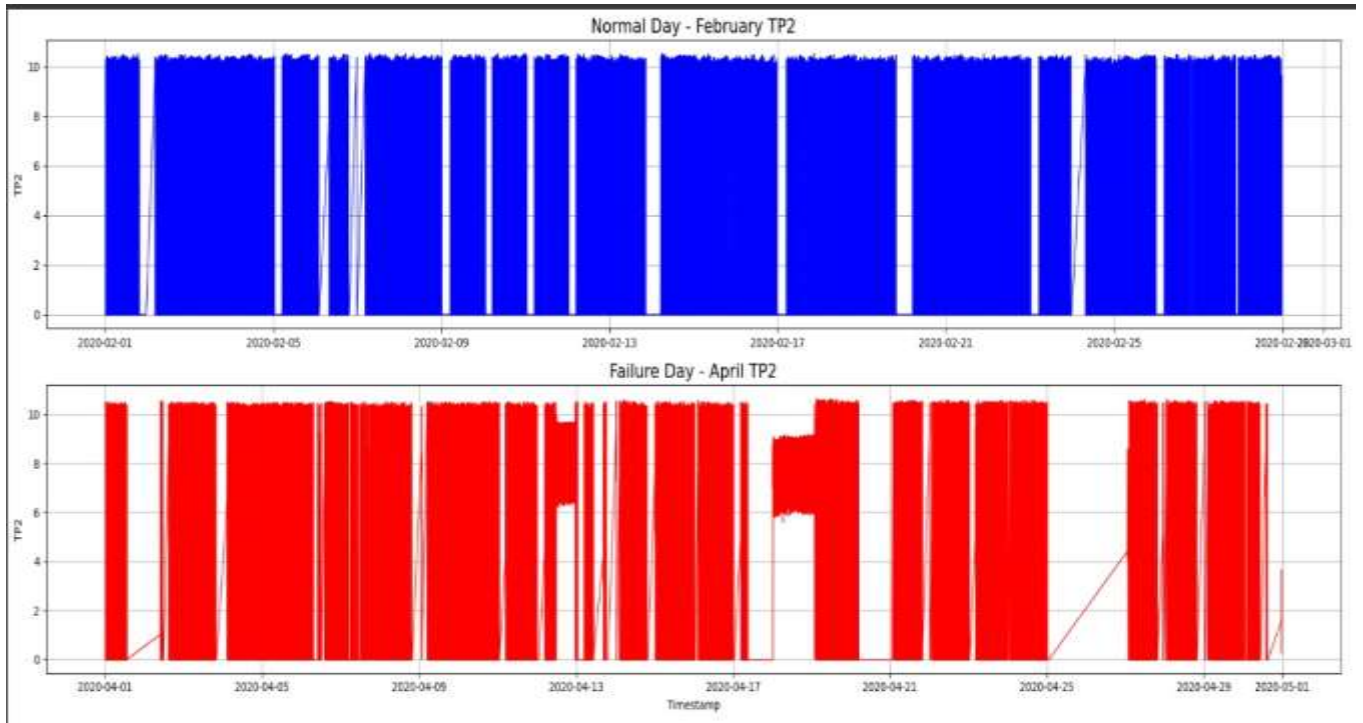**Following are some plots we used to extract some useful information**



*Figure 3: TP2 Analog Signal Comparison: Normal vs. Failure Days*

- Comparison of TP2 analog readings between a normal operating day in February (top, blue) and a failure-affected day in April (bottom, red). The failure day shows significant signal dropouts, irregular fluctuations, and prolonged low values, indicating operational anomalies likely associated with air leak incidents."
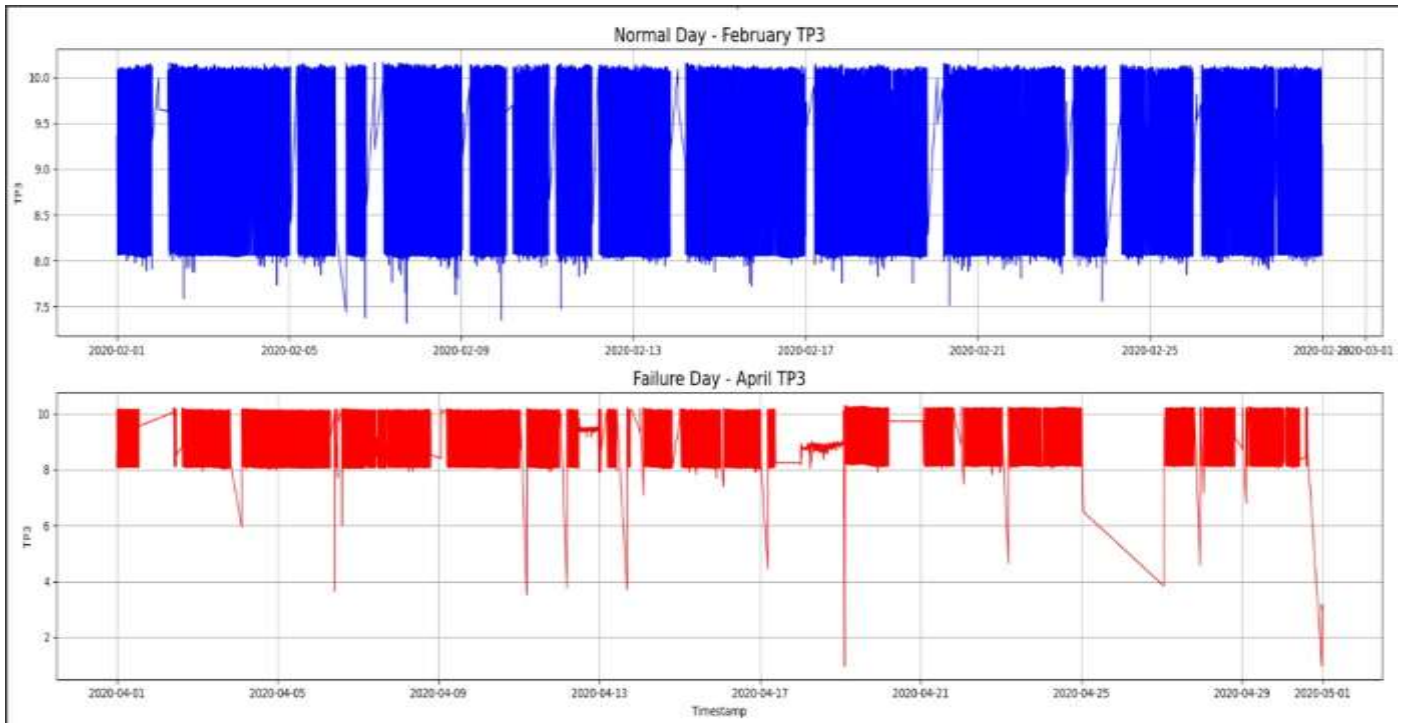
*Figure 4: TP3 Analog Signal Comparison: Normal vs. Failure Days.*

- *The TP3 analog comparison shows stable readings during a normal day in February, with minimal fluctuations. In contrast, the failure day in April exhibits frequent and severe drops, indicating system instability likely linked to operational faults or maintenance events.*
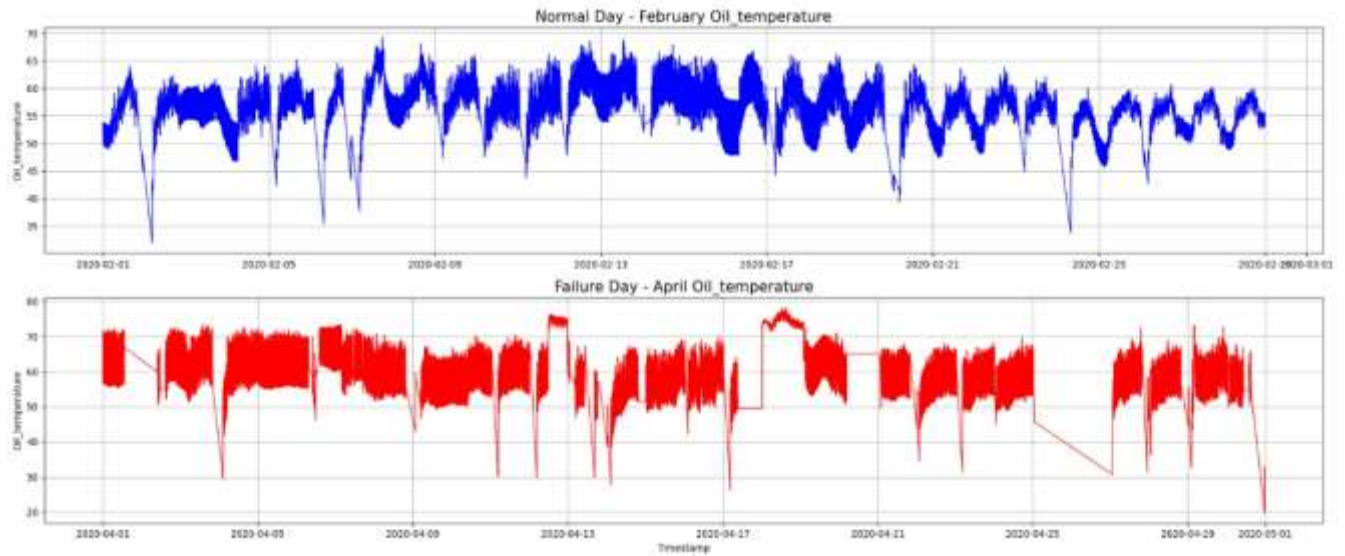
*Figure 5 : Oil-Temperature Analog Signal Comparison: Normal vs. Failure Days.*

- The oil temperature during normal February days remains within a consistent range, showing regular fluctuations. In April's failure period, there are more abrupt drops and irregular spikes, suggesting abnormal operating conditions and possible system faults
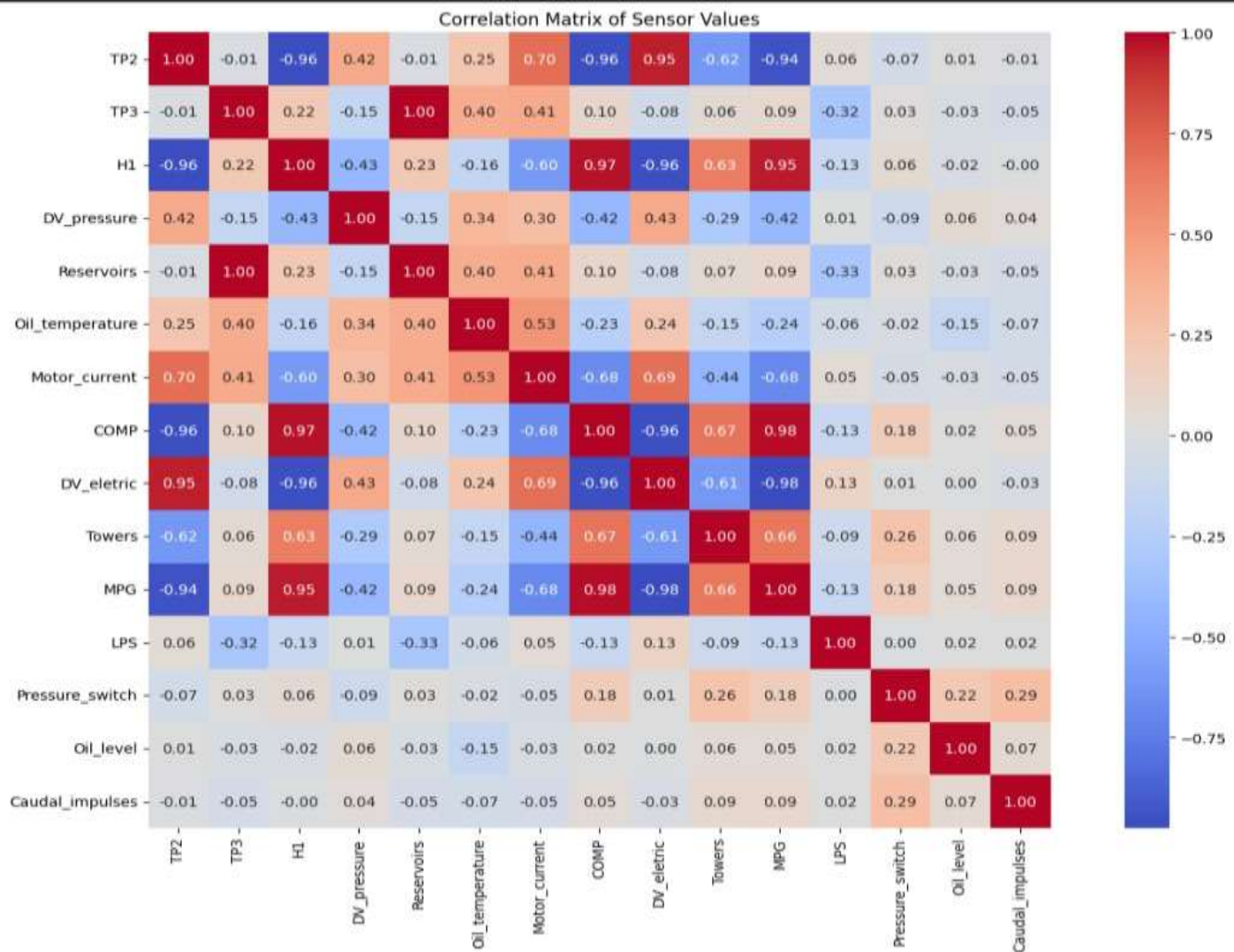
*Figure 6: Correlation Matrix of Sensor Values.*

- *The heatmap shows correlations among different sensor readings. Strong positive correlations (red) indicate sensors with values rising together, while strong negative correlations (blue) indicate inverse relationships. Notably, TP2 and H1, COMP and MPG, and DV_electric and COMP show very strong correlations, suggesting interlinked operational.*

## 2.4 Model Building

1. Train/Test Split

An essential step in building any machine learning model is splitting the dataset into separate portions for training and testing. This ensures that the model's performance is evaluated on data it has never seen before, providing a realistic estimate of how it will perform in production.

Status = 1 → Normal operation (no maintenance required)

Status = 0 → Fault condition (maintenance required)

We divided the dataset chronologically:

Training Data: Earlier months (e.g., March, April, May) were used to train the model.

Testing Data: A later month (e.g., June) was reserved entirely for testing.

2. Procedure:

1. Split the dataset into training and testing sets while preserving time order.

2. Train the model using the training data.

3. Evaluate the trained model on the testing data.

4. Compare predictions with actual labels to measure performance.

## 2.5 Advantages of Train/Test Split

- Ensures unbiased performance evaluation.

- Prevents data leakage, where information from the test set influences training.

- Helps identify overfitting by comparing training and testing accuracy.

- Simulates real-time predictions by keeping future data unseen until testing.

## 2.6 Model Selection

After considering multiple approaches, we selected the Long Short-Term Memory (LSTM) neural network as our primary prediction model. LSTMs are a special type of recurrent neural network (RNN) designed to capture patterns in sequential data, making them well-suited for fault detection in time-series sensor readings.

We also compared the LSTM with simpler models like Logistic Regression and Random Forest to verify its superiority for our use case. The LSTM consistently outperformed these models in capturing temporal dependencies and predicting upcoming faults.

## 2.7 Why LSTM

- Time-Series Capability: It processes data in sequences, learning patterns over time.

- Long-Term Memory: It can remember important trends from earlier readings that might indicate upcoming failures.

- Real-Time Prediction: Works well with streaming data from Kafka and Spark for continuous monitoring.

```
# Train LSTM
epochs = 50
batch_size = 1024
lstm_model.fit(X_train, y_train, epochs=epochs, batch_size=batch_size, verbose=1)

# Save model
lstm_model.save("lstm_model.h5")

Epoch 1/50
627/627 ——————————————— 7s 7ms/step - accuracy: 0.9829 - loss: 0.0852
Epoch 2/50
627/627 ——————————————— 4s 6ms/step - accuracy: 0.9917 - loss: 0.0197
Epoch 3/50
627/627 ——————————————— 3s 5ms/step - accuracy: 0.9972 - loss: 0.0100
Epoch 4/50
627/627 ——————————————— 6s 6ms/step - accuracy: 0.9975 - loss: 0.0090
Epoch 5/50
627/627 ——————————————— 6s 7ms/step - accuracy: 0.9983 - loss: 0.0067
Epoch 6/50
627/627 ——————————————— 4s 5ms/step - accuracy: 0.9984 - loss: 0.0056
Epoch 7/50
627/627 ——————————————— 4s 7ms/step - accuracy: 0.9982 - loss: 0.0059
Epoch 8/50
627/627 ——————————————— 4s 6ms/step - accuracy: 0.9984 - loss: 0.0056
Epoch 9/50
627/627 ——————————————— 3s 5ms/step - accuracy: 0.9985 - loss: 0.0053
Epoch 10/50
```

*Figure 7: LSTM Model Training Performance.*

This figure shows the training log for an LSTM model over 50 epochs with a batch size of 1024.

- The training begins with an accuracy of 98.29% and a loss of 0.0852 in the first epoch.

- By the 2nd epoch, accuracy jumps above 99%, with loss dropping sharply to 0.0197.

- From the 4th epoch onward, accuracy consistently stays around 99.8%, and loss remains extremely low (below 0.006),indicating fast convergence and minimal error.

```
LSTM Metrics:
Accuracy: 0.9195442116542881
Precision: 0.9200676656282643
Recall: 0.999372455306267
F1 Score: 0.9580817572982118
Confusion Matrix:
 [[    21  17294]
  [   125 199064]]
```

*Figure 8: LSTM Model Evaluation Metrics and Confusion Matrix*

This figure presents the performance evaluation of the trained LSTM model:

- Accuracy**:** 91.95% — overall proportion of correct predictions.
- Precision**:** 92.01% — among predicted positives, the percentage that are actual positives.
- Recall**:** 99.94% — the model detects almost all actual positives.
- F1 Score**:** 95.88% — the harmonic mean of precision and recall, indicating balanced performance.

```
for _ in range(1 * 24 * 60 // 60):  # data is 1 row per 10 minutes → 1 days
    X_seq = create_last_sequence(current_df, seq_len)
    y_prob = lstm_model.predict(X_seq, verbose=0)[0][0]
    y_pred = 1 if y_prob > 0.5 else 0  # 1 = no maintenance, 0 = maintenance required

    # Create next timestamp
    next_timestamp = current_df['timestamp'].iloc[-1] + timedelta(minutes=10)

    # Append predicted row (here we fill other features with last known values or dummy)
    next_row = current_df.iloc[-1].copy()
    next_row['timestamp'] = next_timestamp
    next_row['status'] = y_pred
    current_df = pd.concat([current_df, pd.DataFrame([next_row])], ignore_index=True)

    # Store human-readable prediction
    predictions.append({
        "timestamp": next_timestamp,
        "status": "No Maintenance" if y_pred == 1 else "Maintenance Required"
    })

# Convert to DataFrame
pred_df = pd.DataFrame(predictions)
print(pred_df)
```
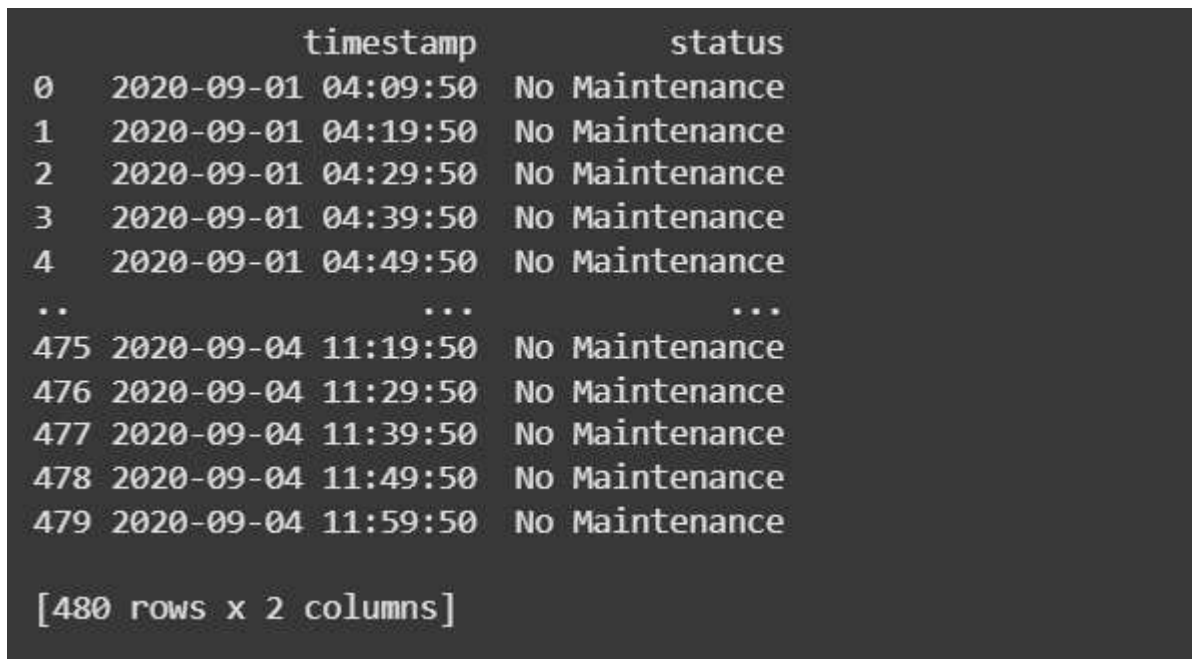
*Figure 9: LSTM-Based Future Maintenance Prediction*

- This code predicts maintenance requirements for the next 24 hours in 10-minute intervals using a trained LSTM model. At each step, it generates a new timestamp, predicts the status (No Maintenance or Maintenance Required), and appends it to the dataset. Predictions are stored in a human-readable format and compiled into a DataFrame for analysis. This enables proactive, step-by-step forecasting for maintenance planning.

```
           timestamp              status
0    2020-09-01 04:09:50    No Maintenance
1    2020-09-01 04:19:50    No Maintenance
2    2020-09-01 04:29:50    No Maintenance
3    2020-09-01 04:39:50    No Maintenance
4    2020-09-01 04:49:50    No Maintenance
..                  ...               ...
475  2020-09-04 11:19:50    No Maintenance
476  2020-09-04 11:29:50    No Maintenance
477  2020-09-04 11:39:50    No Maintenance
478  2020-09-04 11:49:50    No Maintenance
479  2020-09-04 11:59:50    No Maintenance

[480 rows x 2 columns]
```

*Figure 10:  Metro Failure Status Monitoring*

- This figure illustrates a time-stamped dataset used for metro failure detection. Each row captures the status of the system at a specific moment, indicating whether the metro required maintenance or was operating normally. This continuous record of "No Maintenance" signals helps monitor trends and spot irregularities. The structure supports machine learning approaches by providing sequential operational data for model training and prediction. Such tabular time series is foundational for automated failure detection in metro systems.

## 2. Requirements Specification

### 2.2 Hardware Requirement:

- 500 GB hard drive (Minimum requirement)

- 8 GB RAM (Minimum requirement)

- PC x64-bit CPU
- PC x64-bit CPU PC with x64-bit CPU (multi-core preferred for parallel processing)

- GPU (Optional but recommended for faster model training, e.g., T4 GPU)

### 2.3 Software Requirement:

- Windows/Mac/Linux

- Python-3.9+

- Jupyter Notebook or VS Code

- Apache Kafka (for data streming)

- Apache Spark with PySpark (for big data processing)

## 1.7  Required Libraries:

- Numpy

- Pandas

- Matplotlib

- Scikit-learn

- TensorFlow/Keras.

- PySpark

- Kafka-Python

19

### 3. Conclusion

- This project successfully demonstrated the application of machine learning techniques, particularly LSTM-based models, to predict failures in metro air compressor systems using sensor time-series data.
- By preprocessing real-time streaming data through Kafka and PySpark , and deploying trained models for anomaly detection and failure forecasting, the system provides timely maintenance alerts that can significantly reduce downtime and maintenance costs.
- The results show promising accuracy and reliability, validating the effectiveness of the ML approach in predictive maintenance scenarios. Overall, this project lays a strong foundation for scalable, data-driven maintenance solutions in industrial systems.
- .

### 4. Future Scope

• Online and Incremental Learning: Implement online learning techniques to update the model continuously as new data arrives, ensuring the model adapts to changing system behavior without full retraining.

• Anomaly Detection Enhancement: Combine unsupervised anomaly detection methods with supervised models for better identification of rare or unseen failure modes.

• Model Deployment Optimization: Optimize the model for edge deployment, reducing size and computational requirements to enable real-time inference directly on devices near the metro components.

• Advanced Model Architectures: Experiment with more sophisticated architectures such as Transformer-based models, Graph Neural Networks (GNNs), or hybrid models combining LSTM with CNNs to better capture complex temporal patterns

## 5. References

- MetroPT3 Dataset – Predictive Maintenance Data for Metro Trains. Available at: https://www.kaggle.com/datasets
- Apache Kafka Documentation – Real-time Data Streaming. Available at: https://kafka.apache.org/documentation/
- Apache Spark Documentation – Structured Streaming. Available at: https://spak.apache.org/structured-streaming/
- TensorFlow Documentation – Deep Learning Framework. Available at: https://www.tensorflow.org/api_docs
- Keras Documentation – High-level Neural Networks API. Available at: https://keras.io/api/