

Metro PT3 Predictive Maintenance

Fundamentals of Data mining: Dr Seif Eldawlatly

Malak Gaballa – Masa Tantawy – Moustafa El Mahdy



Table of Contents

01 Introduction

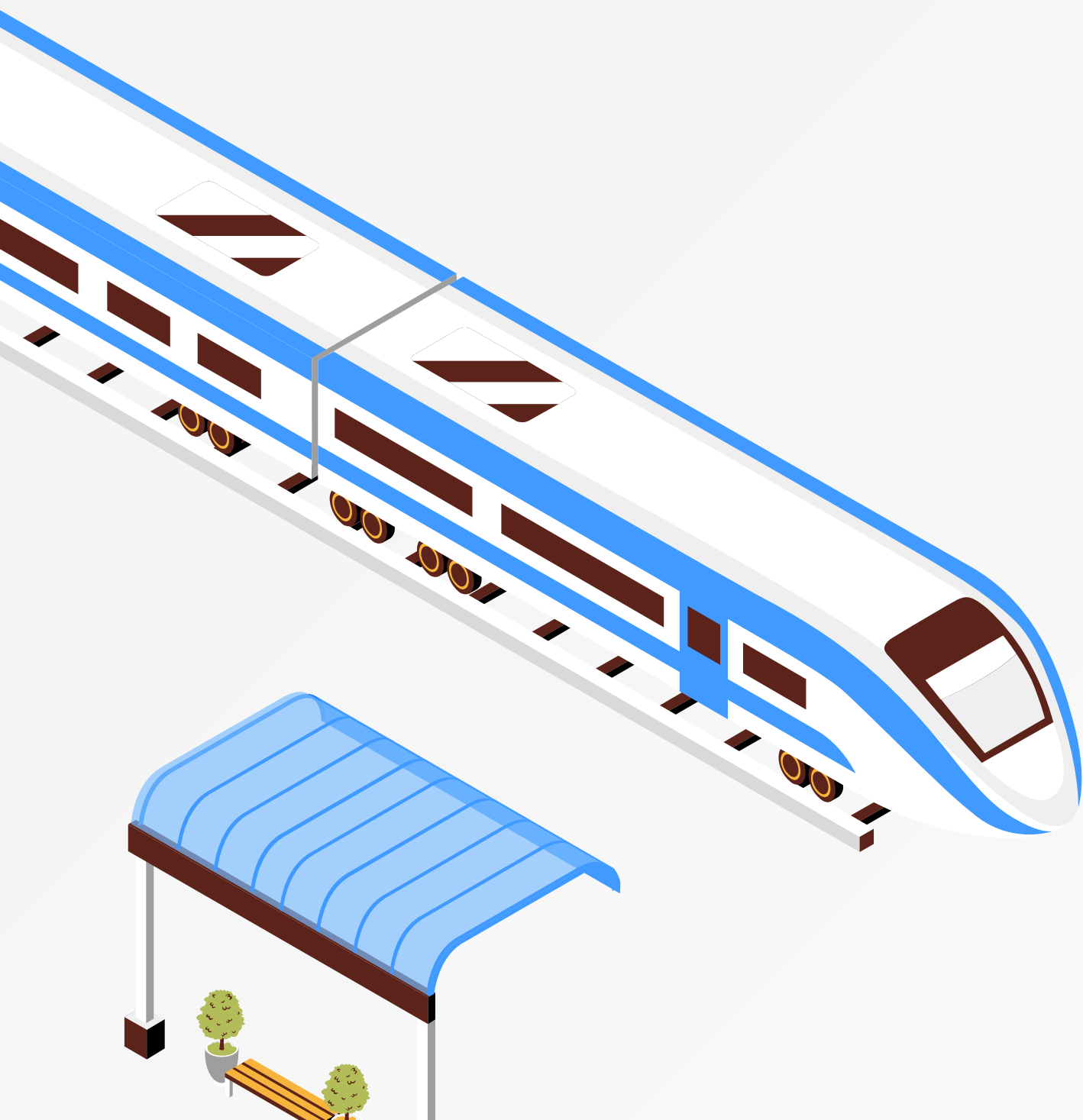
Data Description and
Preprocessing

02 Methodology and Results

Applied Approaches,
Attempted Techniques &
Evaluation

03 Conclusion

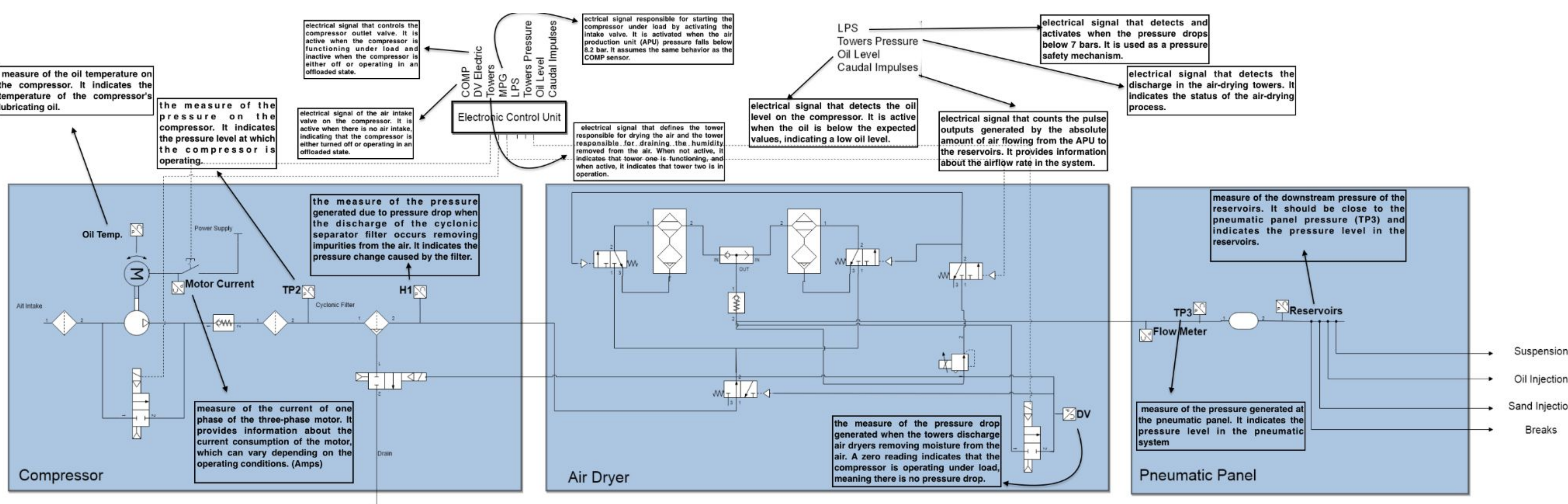
Final
Recommendations



The Dataset

- Metro Train Data -

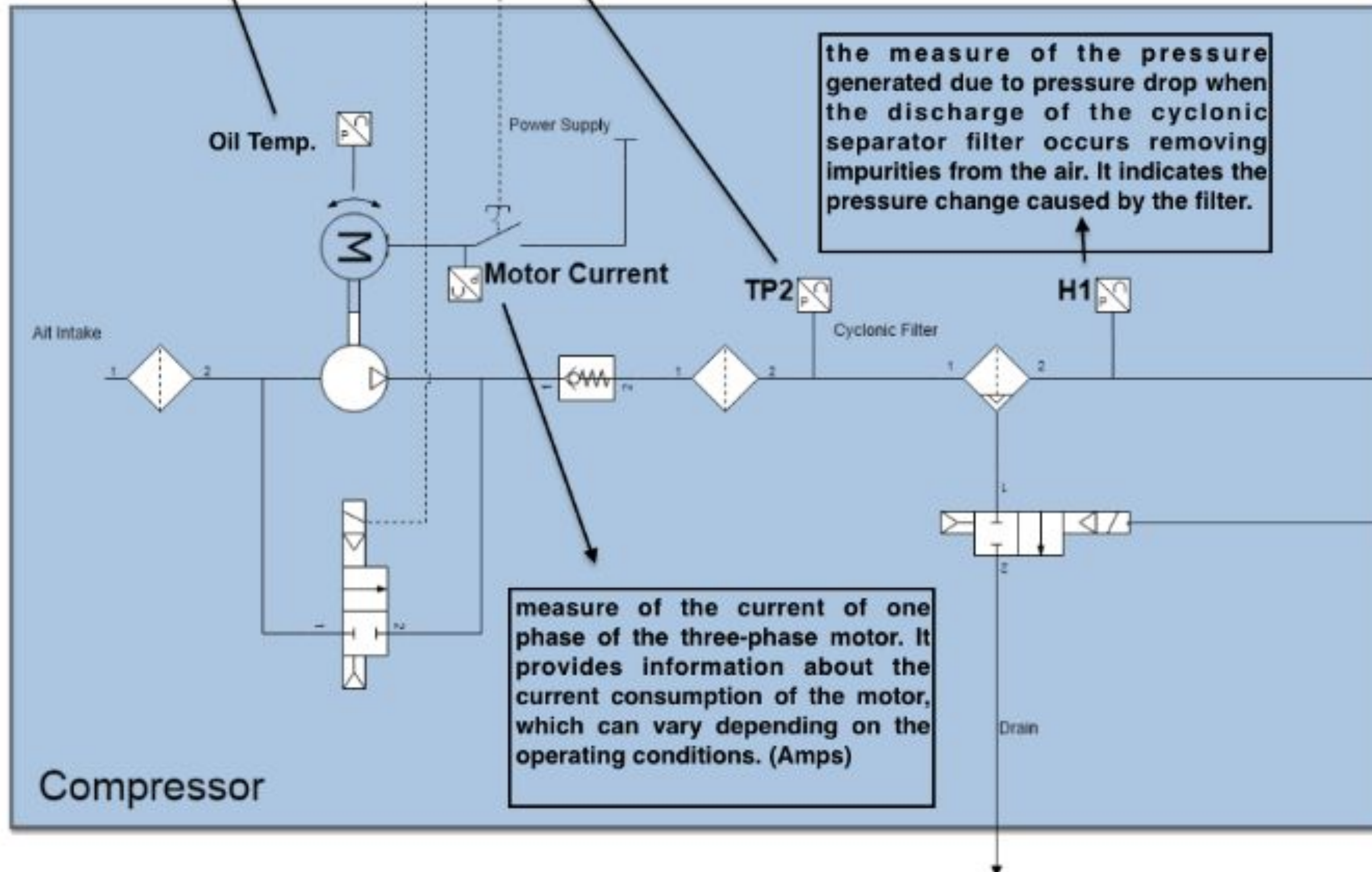
- It consists of Compressor's Air Product Unit(APU) readings.
- Collected at 1Hz from February to August 2020
- Pressure, temperature, motor current consumption, and air intake valves
- 1,516,948 readings or instances and 16 Features
 - Timestamp
 - 7 Analogue Sensors (Continuous)
 - TP2 , TP3, HI , DV Pressure , Reservoirs , Motor Current , Oil Temperature
 - 8 Digital Sensors (Binary)
 - Comp , DV Electric , Towers , MPG , LPS , Pressure Switch , Oil Level , Caudal Impulses
- Unlabelled data (Initially)
 - Aim: detect observations with APU failure (2.05% of the data)

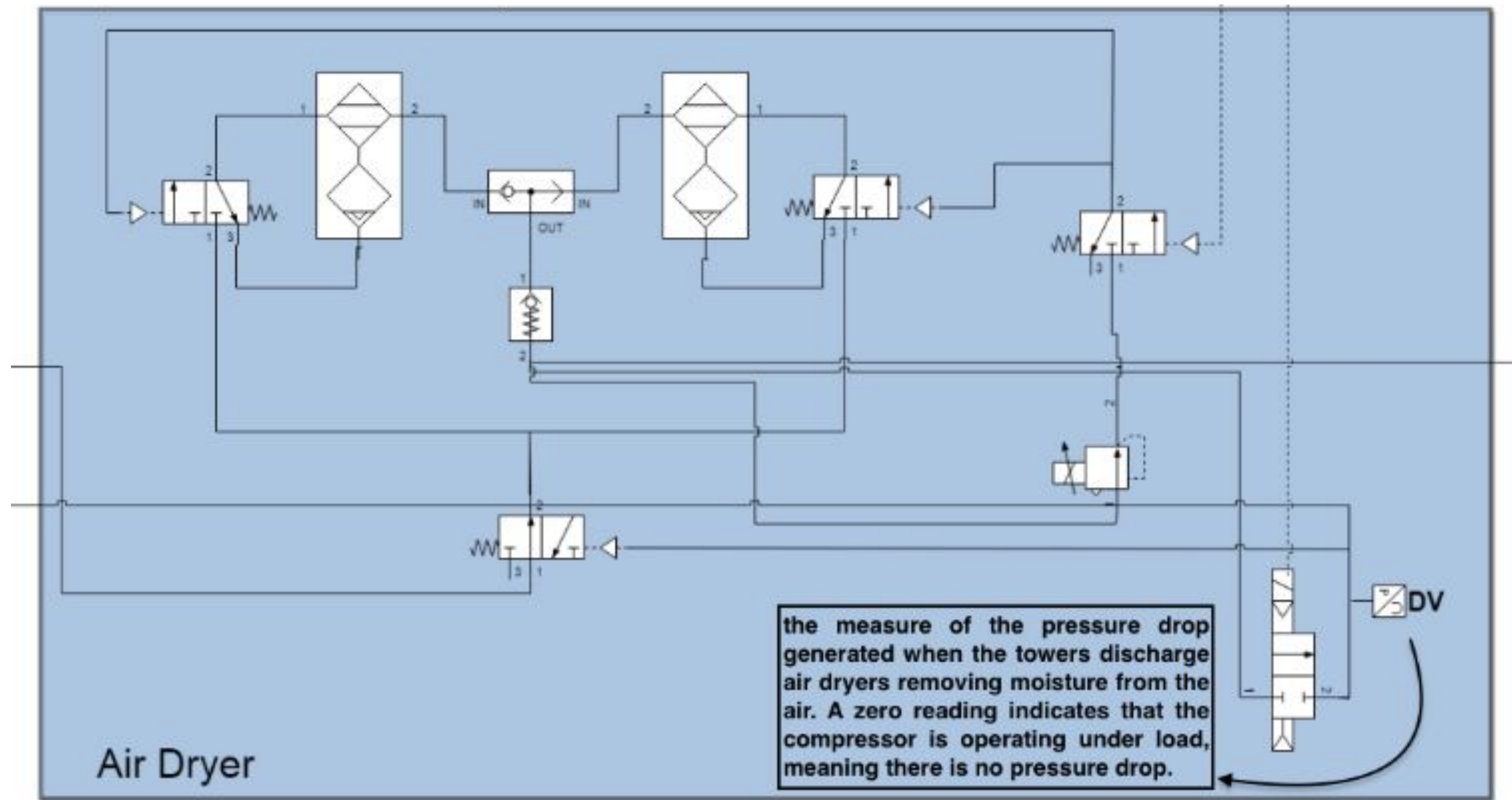


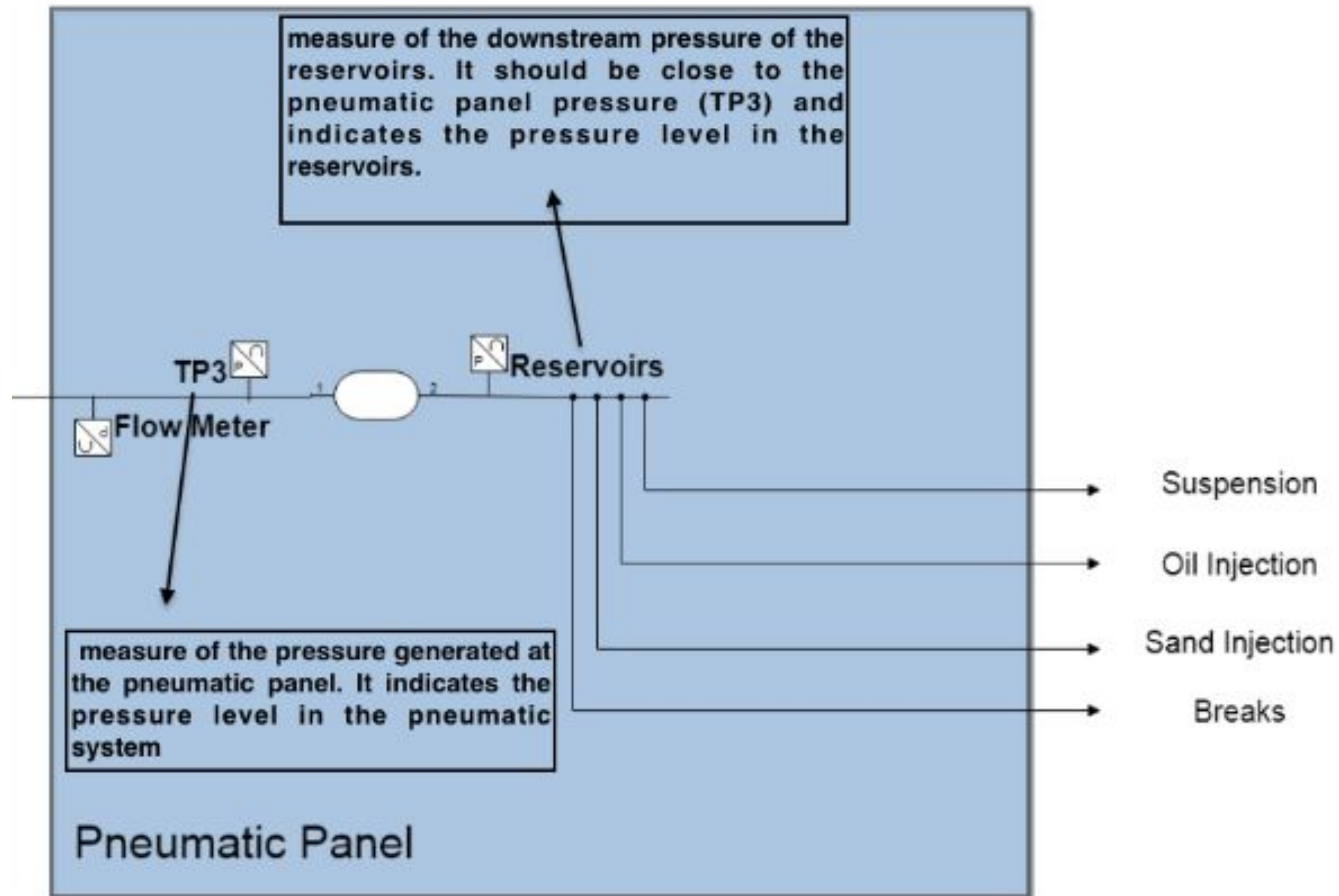
measure of the oil temperature on the compressor. It indicates the temperature of the compressor's lubricating oil.

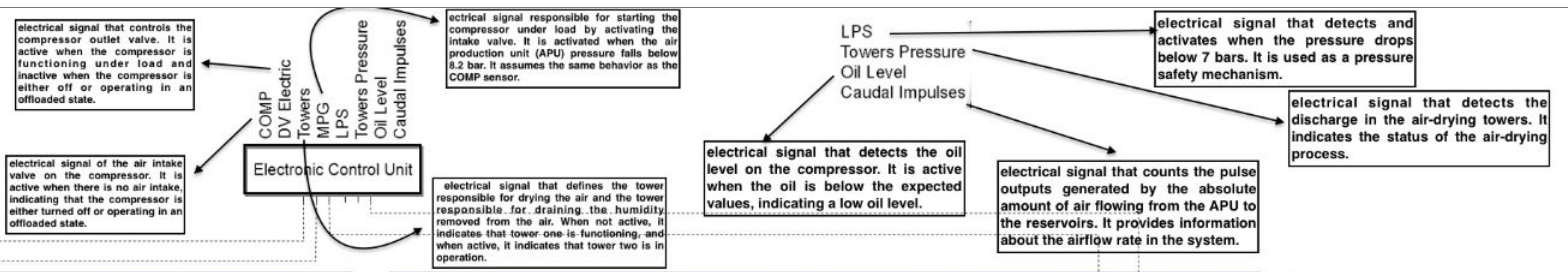
the measure of the pressure on the compressor. It indicates the pressure level at which the compressor is operating.

the measure of the pressure generated due to pressure drop when the discharge of the cyclonic separator filter occurs removing impurities from the air. It indicates the pressure change caused by the filter.









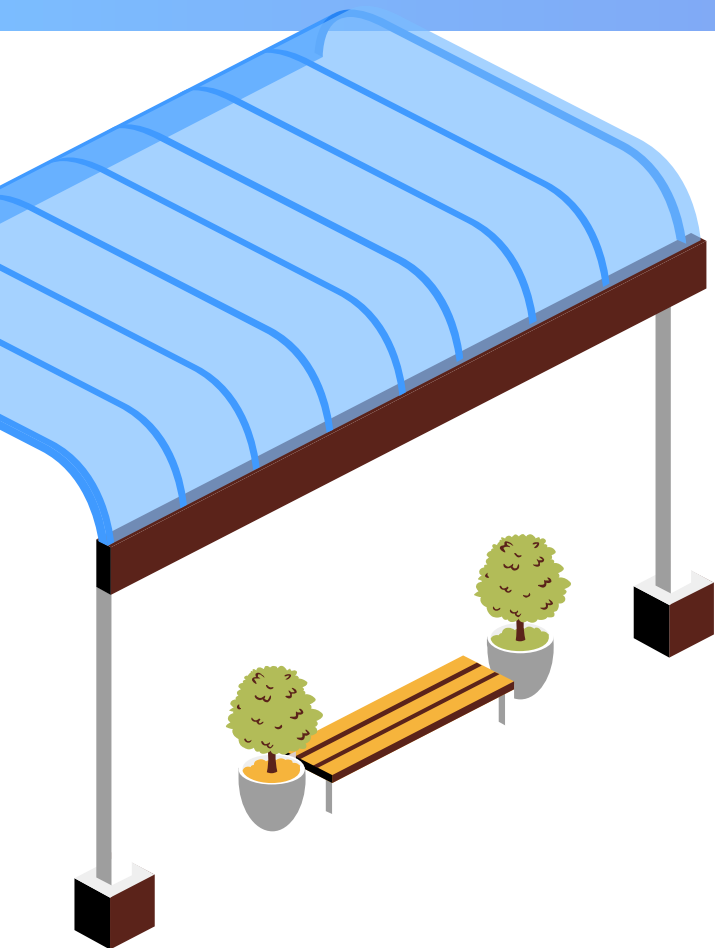
The Dataset

- Data Labelling -

- Aim: detect observations with APU failure
- Data is labelled based on air leaks from maintenance reports
- **Label 'Airleak':** 0 → no airleak, 1 → airleak (flagged)
- 29,865 readings (2.05% of the data)

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

Methodology



PCA on Full data
2 PCA components
1,459,475 Obs.
2.05% anomalies

Approach 2

Balanced data:
All minority + RUS majority
(59,730 , 15)
50% anomalies

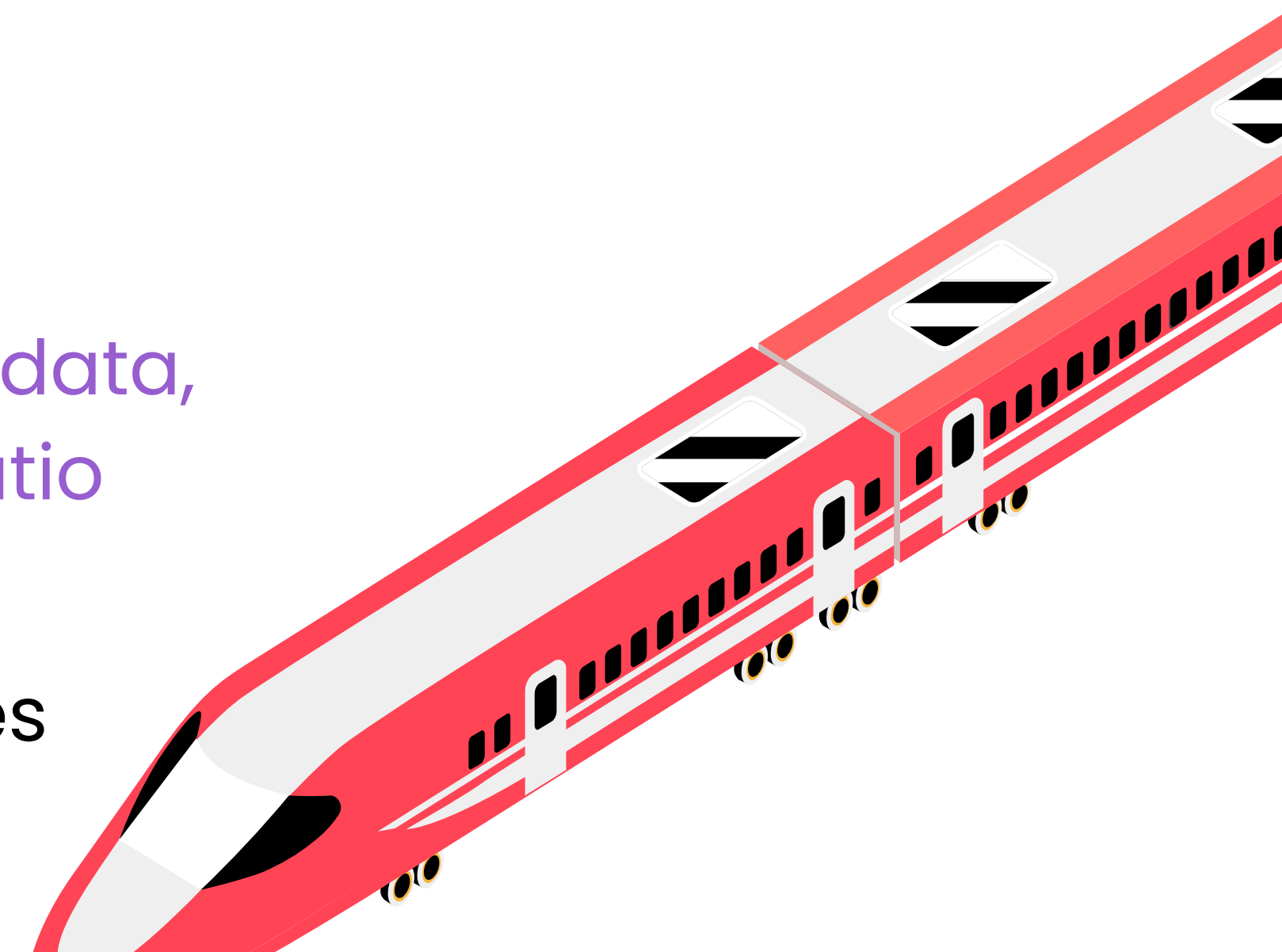
Approach 4

Approach 1

Full data
(1,459,475 , 15)
2.05% anomalies

Approach 3

Half the data:
random sample of data,
preserving the ratio
(729,738 , 15)
2.07% anomalies



Methodology

- Anomaly detection - (Approaches 1,2,3)

1. Clustering-based Approach

- k-Means clustering (euclidean distance)
- Hierarchical Clustering (ward, single, complete, ...)
- Spectral Clustering

2. Distance-based Approach

- Fifth Nearest Neighbour (euclidean distance)
- Mahalanobis Distance
- BACON

3. Density-based Approach

- DBSCAN: Density Based Spatial Clustering of Applications with Noise

Methodology

- Classification - (Approach 4)

1. Clustering-based Approach

- k-Means clustering (euclidean distance)
- Hierarchical Clustering (ward, single, complete, ...)
- Spectral Clustering

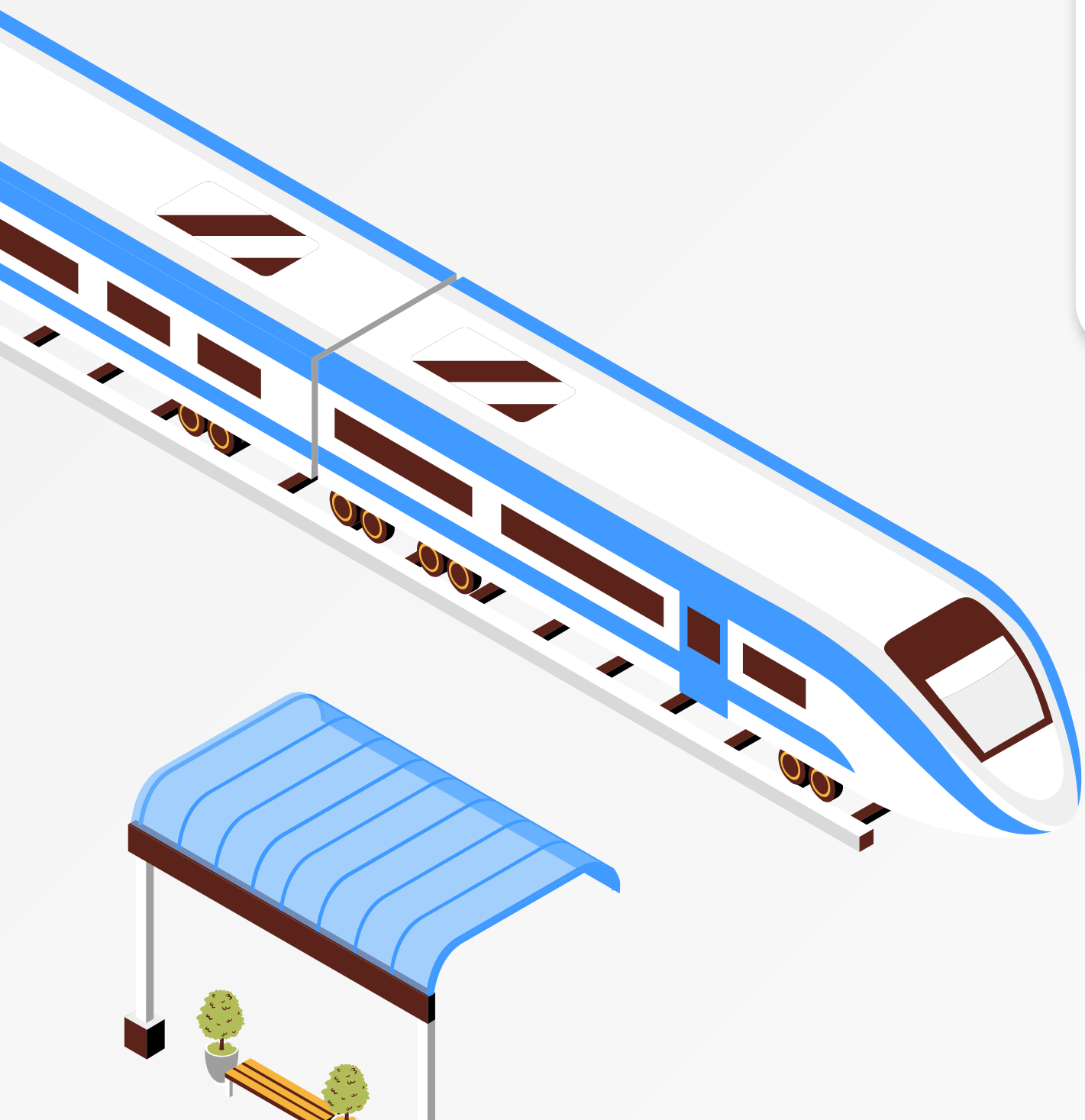
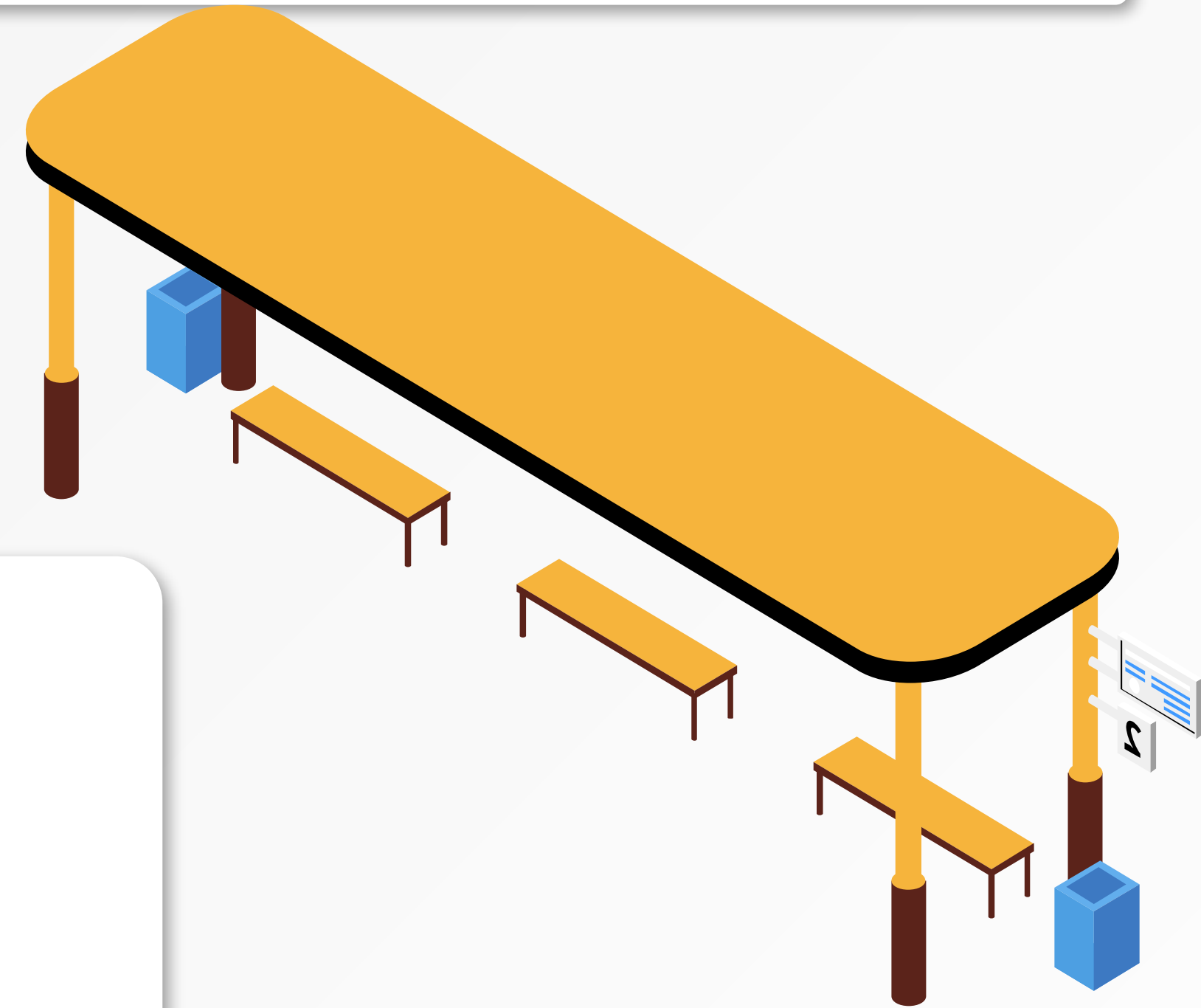
2. Machine Learning Classifiers

- Random Forest
- K-Nearest Neighbour (K-NN)
- Logistic Regression
- Support Vector Machines (SVM)
- Least squares
- Perceptron

Methodology

Performance Metrics

- Precision
(true positive rate)
- Recall
(false positive rate)
- F1-Score
(weighted average of both)



Approach 1: Full data

1. Clustering

	Precision	Recall	F1-score
<u>K-means</u>	0.05	1.00	0.09
<u>Hierarchical Clustering</u>	Crashed		
<u>Spectral Decomposition</u>	Crashed		

2. Distance-based Approach for Anomaly Detection

	Precision	Recall	F1-score
<u>Fifth Nearest Neighbour</u>	0.09	0.18	0.12
<u>Mahalanobis Distance</u>	0.21	1.0	0.35
<u>BACON</u>	0.07	1.00	0.14

3. Density-based Approach for Anomaly Detection

DBSCAN (Density Based Spatial Clustering of Applications with Noise)

	Precision	Recall	F1-score
eps= 0.2, min_samples= 4	Crashed		
eps= 0.2, min_samples= 3	Crashed		

Approach 2: PCA (2 components) on full data

1. Clustering

	Precision	Recall	F1-score
<u>K-means</u>	0.13	1.00	0.24
<u>Hierarchical Clustering</u>	Crashed		
<u>Spectral Decomposition</u>	Crashed		

2. Distance-based Approach for Anomaly Detection

	Precision	Recall	F1-score
<u>Fifth Nearest Neighbour</u>	0.13	0.17	0.14
<u>Mahalanobis Distance</u>	0.20	1.00	0.33
<u>BACON</u>	0.06	1.00	0.11

3. Density-based Approach for Anomaly Detection

DBSCAN (Density Based Spatial Clustering of Applications with Noise)

	Precision	Recall	F1-score
eps= 0.2, min_samples= 4	Crashed		

Approach 3: Random Sample (1/2 the data, same ratio)

1. Clustering

	Precision	Recall	F1-score
<u>K-means</u>	0.05	1.00	0.09
<u>Hierarchical Clustering</u>	Crashed		
<u>Spectral Decomposition</u>	Crashed		

2. Distance-based Approach for Anomaly Detection

	Precision	Recall	F1-score
<u>Fifth Nearest Neighbour</u>	Crashed		
<u>Mahalanobis Distance</u>	0.22	1.00	0.35
<u>BACON</u>	0.07	1.00	0.13

3. Density-based Approach for Anomaly Detection

DBSCAN (Density Based Spatial Clustering of Applications with Noise)

	Precision	Recall	F1-score
eps= 0.2, min_samples= 4	Crashed		
eps= 0.2, min_samples= 3	0.04	0.05	0.05

Approach 4: Balanced Data

1. Clustering

	Precision	Recall	F1-score
<u>K-means</u>	0.92	1.00	0.96
<u>Hierarchical Clustering</u>	Crashed		
<u>Spectral Decomposition</u>	Crashed		

2. Machine Learning Classifiers

Results are based on cross-validation on 5-fold

	Precision	Recall	F1-score
<u>Random Forest</u>	1.00	0.93	0.96
<u>K-NN (3)</u>	0.99	0.93	0.96
<u>Logistic Regression</u>	0.98	0.94	0.96
<u>SVM</u>	0.98	1.00	0.99
<u>Least Squares</u>	0.98	0.94	0.96
<u>Perceptron</u>	0.98	0.98	0.98

Conclusion

Finally, it can be concluded that **machine learning classifiers performed the best.**

For future recommendations, it would be valuable to attempt techniques that crashed in our application, or attempt dimensionality reduction (PCA) on the sampled datasets.