

**COMPARATIVE ANALYSIS OF AI AND XAI TECHNIQUES
ACROSS DIVERSE PERFORMANCE METRICS ON VARIED
DATASETS**

**Prashant C. Dhas¹, Dr. Vidula V. Meshram², Dr. Parikshit
N. Mahalle³**

¹Ph.D. Research Scholar, Vishwakarma Institute of Information
Technology, Pune-48,

Savitribai Phule Pune University (SPPU), Pune

² Research Supervisor, Assistant Professor, Department of Computer
Engineering,

Vishwakarma Institute of Technology, Pune-48, Savitribai Phule Pune
University (SPPU), Pune

³ Professor, Department of Artificial Intelligence and Data Science,
Vishwakarma Institute of Technology, Pune-48, Savitribai Phule Pune
University (SPPU), Pune

prashant.221p0054@viit.ac.in¹, vidula.meshram@vit.edu²,
aalborg.pnm@gmail.com³

Abstract

This study paper compares Artificial Intelligence (AI) and Explainable Artificial Intelligence (XAI) methods by looking at how well they work in different areas using the "Machine Failure Prediction using Sensor Data" dataset. The goal is to find out how different models work and how XAI can make models easier to understand and more open. The dataset used has both features that have been preprocessed using standard scalar normalisation and data that has been balanced using SMOTE. This makes sure that the model training is strong and fair. For this study, a lot of testing was done with different machine learning classifiers, such as the Gradient Boosting Classifier, the Isolation Forest, the Support Vector Machine (SVM), the Optimised Isolation Forest (OIF), and mixed models like the 1D CNN-OIF and CNN-LSTM. Standard performance measures were used to judge each model. These were based on exploratory data analysis (EDA) and a feature selection process that looked for features that had an association with the goal label of more than 0.3. XAI methods, such as SHAP (SHapley Additive Explanations) and LIME (Local Interpretable Model-agnostic Explanations), were used to understand the AI models, giving details about the features that went into making the models and the choices they made. Accuracy and loss curves, confusion vectors, and a comparison of model outcomes were used to look at the data. The combined CNN-LSTM model did better than the others in this study in terms of accuracy and ease of use, making it the chosen model

for further research. Thanks to the use of XAI methods, this model not only did a better job of making predictions, but it also gave us a lot of information about how decisions are made.

Keywords: Artificial Intelligence, Explainable AI, Machine Learning, Performance Metrics, Hybrid Models, Sensor Data Analysis

1. Introduction

Artificial Intelligence (AI) has become a key technology that is advancing many fields, such as industry, healthcare, and banking. It is a key part of the growing fields of data science and machine learning. In the field of machine learning fashions, AI structures have shown that they could manage challenging jobs like making selections in actual time and planning preventative preservation. Despite the fact that those systems paintings, they regularly work like "black packing containers," wherein the selection-making techniques aren't clear. This would make belief and duty troubles more likely. This requires the use of Explainable synthetic talent (XAI) methods, which are supposed to help people understand what AI fashions do. XAI is in particular important in conditions in which safety and dependability are essential, like in predictive preservation fashions used in enterprise. These models can tell whilst device will wreck down and layout restore for the right time, which maintains things running easily and makes the equipment ultimate longer [1]. The "machine Failure Prediction the usage of Sensor information" dataset from Kaggle was used in this have a look at. It has a number of sensor reviews which are useful for predicting while a device will smash down. This dataset has quite a few distinctive sensor outputs which can be had to construct sturdy machine learning models that may correctly predict failures [2].

The primary goal of this take a look at is to compare and contrast extraordinary AI and XAI strategies on the way to discover which ones paintings fine based totally on one-of-a-kind performance measures. A lot of complex learning strategies is used in this take a look at. Those consist of Gradient Boosting Classifier, Isolation forest, support vector machine (SVM), and Optimised Isolation forest (OIF). Also used are new blended models like 1D CNN-OIF and CNN-LSTM. Our desire of these models was primarily based on how properly they could cope with the complex and high-dimensional sensor information. XAI strategies like SHAP (SHapley Additive factors) and LIME (nearby Interpretable model-agnostic motives) are used in this have a look at to solve the problem of interpretability. These strategies are meant to disclose how every feature impacts the version's preference, which makes AI structures more open and dependable. As an instance, SHAP values supply a whole lot of data approximately the prediction options, which allows humans recognize and believe the consequences of the predictive models. Further, LIME uses locally replacement fashions to explain any classifier's effects in a technique that is clear and correct [3].

An organised way to model review is used in this study's methodology. The dataset goes through a strict preparation phase that gets rid of duplicates, normalises the data, and balances it using SMOTE to make training and testing the model fairer. An exploratory data analysis (EDA) is done to find traits that have a strong relationship with the goal variable. This makes

the training process more focused and effective. The success of these models is carefully evaluated using a number of different measures. Some of these are F1 score, accuracy, sharpness, and memory. Both how well each model can predict the future and how well it can be analysed using XAI methods are used to rate its usefulness. This two-part review gives a full picture of how well each model does by looking at both how well it completes tasks and how clear it is when describing choices. A lot of useful information was learnt about how different AI and XAI setups work through this study. From the results, it looks like mixed models, especially the CNN-LSTM, did better in terms of accuracy and loss measures. More importantly, these models got a lot better when XAI methods were used on them. These techniques made them much easier to understand and more reliable.

The major contribution of paper is given as:

- Enhanced Model Understanding through XAI: Applied SHAP and LIME to demystify AI decisions, improving transparency and trust in predictive maintenance models.
- Benchmarking Hybrid AI Models: Evaluated and benchmarked hybrid models like 1D CNN-OIF and CNN-LSTM, showcasing superior accuracy and better integration with XAI.
- Comprehensive Performance Metrics Analysis: Detailed analysis of performance metrics to identify top-performing models, crucial for effective real-world applications in machine failure prediction.

2. Related Work

A lot of people in academia and business are interested in how to use AI and XAI methods together in different areas. People are interested in this because we need to make AI systems more open and responsible, especially in areas where safety and dependability are very important. The ability of AI to see possible problems before they happen can greatly cut down on downtime and operational costs in predictive maintenance [3]. This makes the technology essential for modern manufacturing processes. New developments in machine learning algorithms have made prediction models more accurate. Gradient Boosting Classifiers and Support Vector Machines (SVMs), for example, are widely used because they are good at dealing with big and complicated datasets [4]. People really like these models because they can accurately predict what will happen and are reliable, which is very important in the high-stakes world of machine failure forecast. Isolation Forests are also used for anomaly detection, which is another step forward in technology because they help us understand strange trends that could mean that something is about to break [5].

Some combination models, like the CNN-LSTM and 1D CNN-OIF, help make predictions greater accurate. This suggests that AI is constantly converting and improving. Studies have shown that these fashions now not only make predictions greater accurate, however they also help us understand the temporal and spatial patterns of datasets better [6]. the use of these models in sensor information analysis to are expecting machine failure has proven that they are a good deal greater correct and use loads less computing power than conventional gadget

mastering techniques [7]. Explainable AI (XAI) has emerged as an essential region of have a look at because it solves the hassle of "black container" AI systems. Inside the past few years, loads of studies has been finished on techniques like SHAP and LIME, which show how extraordinary traits can affect version effects. These techniques are very essential for making sure that AI fashions are truthful and dependable and for constructing faith among customers. SHAP, particularly, offers regular and accurate solutions, which is why it's far the pleasant choice for research that need to do thorough feature evaluation [8, 9]. numerous research have looked at the differences and similarities among AI and XAI methods, especially when it comes to machine learning models and the way clean they are to understand. A variety of the time, those studies have a look at the change-offs among how challenging it's far to recognize a model and how complicated it's miles. As an instance, simpler fashions like decision bushes can be explained by using basic XAI techniques, but greater complex models like deep neural networks need greater advanced XAI strategies to provide an explanation for how they make choices [10, 11]. This brings up the cutting-edge discussion within the AI discipline approximately the high-quality approaches to balance speed and openness.

A lot of study has also been made about how AI and XAI can be used in sensor-based data systems. As sensor data is multidimensional and contains noise, it presents special problems for AI models to solve. Preprocessing steps like data normalisation and feature selection are very important for getting this kind of data ready for training and analysing models [12, 13]. There is a lot of evidence that SMOTE can balance datasets to help machine learning models do better on data that isn't balanced [14, 15]. Aside from the technical elements, AI and XAI are also getting more attention for their legal and moral effects. Due to the need for models to follow legal and moral rules, there has been a lot of study into making rules and frameworks for the responsible use of AI technologies. Quite a few the times, those studies strain how important it is for AI structures to be open, accountable, and sincere, especially in important regions like healthcare and public safety [16, 17].

Overall performance measures like accuracy, precision, memory, and F1 rankings are regularly used to choose AI and XAI strategies. Those ratings supply a full photograph of ways nicely the models paintings. This approach now not only allows find the pleasant models; however it also allows find places in which adjustments are wished. Latest research indicates that combining performance evaluation with XAI motives could make AI structures more solid and dependable, given that those insights make it easier to preserve improving and refining the models [18, 19]. Sooner or later, one routine topic in modern look at is how AI and XAI trade social values and commercial enterprise strategies. Researchers have regarded into how these technologies can alternate industries by making it less difficult to make selections and going for walks corporations more effectively. AI has a big potential to force innovation, but its wider outcomes must additionally be carefully notion thru, which includes the chance of sudden results and the need for controls to stop misuse [20, 21].

Table 1: Summary of the related work to the study of AI and XAI techniques

Study Focus	Model/Technique	Key Findings	Performance Metrics	Data Type
Predictive Maintenance	General AI	AI reduces downtime and operational costs	Not specified	Manufacturing sensor data
Machine Learning Advancements	Gradient Boosting, SVM	Effective in large, complex datasets	Predictive accuracy	Various datasets
Anomaly Detection	Isolation Forest	Provides insights into unusual patterns indicating failures	Robustness	Manufacturing sensor data
Hybrid AI Models	CNN-LSTM, 1D CNN-OIF	Enhance predictive performance and understanding of data features	Accuracy, computational efficiency	Sensor data
Explainable AI Techniques	SHAP, LIME	Facilitate understanding of individual feature contributions	Consistency and accuracy	Varied
Model Complexity	Deep neural networks	Require sophisticated XAI for transparency	Trade-off analysis	Varied
Data Preprocessing	SMOTE	Critical for preparing high-dimensional, noisy data for effective training	Improved model performance	High-dimensional sensor data
Regulatory and Ethical Aspects	General AI and XAI	Emphasis on transparency, accountability, user trust	Not specified	Critical applications
Evaluation of AI/XAI	General AI and XAI	Comprehensive assessment through standard metrics to identify effective models	Accuracy, precision, recall, F1	Varied

Societal Impact	General AI and XAI	AI transforms industries, enhances decision-making but requires safeguards	Not specified	Varied
-----------------	--------------------	--	---------------	--------

3. Methodology

A. Dataset Used

The "Machine Failure Prediction using Sensor statistics" dataset on Kaggle is a whole set of sensor reports that are meant to assist expect when machines will smash down in an industrial putting. These records may be very important for making predictive maintenance models, which assist predict breakdowns before they take place, which lowers the cost of upkeep and downtime. Some of the sensor outputs that make up the facts are temperature, stress, vibration, spinning, and sound emissions, among other bodily and operational elements. Every item within the record is a picture of sensor values at a certain time, together with a note announcing whether there was a loss or no longer.

	footfall	tempMode	AQ	USS	CS	VOC	RP	IP	Temperature	fail
0	0	7	7	1	6	6	36	3	1	1
1	190	1	3	3	5	1	20	4	1	0
2	31	7	2	2	6	1	24	6	1	0
3	83	4	3	4	5	1	28	6	1	0
4	640	7	5	6	4	0	68	6	1	0

Figure 1. Dataset Sample

This simple binary classification makes it easy to use machine learning models to learn from the trends that show up before something fails. The sensors give us a lot of different features, so we need to be very careful with the preparation to make sure the model training is strong. For better results, it's important to use techniques like normalisation and feature selection to deal with the richness of the data. Because it is complicated and useful in the real world, the dataset is a great way to test how well different AI and XAI methods work. Models find it hard to understand and learn from sensor data's complicated trends. This makes it a useful tool for studying predictive maintenance in the manufacturing industry. The figure 1 demonstrate of data comes from different monitor reports that are related to how the machine works and are probably used for predicted maintenance. Different types of columns, like "footfall,"

"tempMode," "AQ" (Air Quality), "USS" (Ultrasonic Sensor Signal), "CS" (Current Sensor), "VOC" (Volatile Organic Compounds), "RP" (Rotation Per Minute), and "IP" (Impact Pressure), give information in more than one direction. The last two columns show "Temperature" and "fail," with "fail" showing whether a loss happened (1) or not (0). Each row shows a separate measurement event. This data format makes it possible to do complicated analysis to predict when a machine will break down based on different monitor inputs, which is very important for keeping operations running smoothly.

B. Data Pre-Processing

a. Drop Duplicates: This is an important part of preparing the data, especially for datasets with a lot of repeated rows that could throw off the results of later studies. In predictive maintenance, each data point should ideally show a different case of how the machine is working. Duplicate records can cause false trends and mistakes in training the model. Getting rid of these copies makes sure that the machine learning models aren't biased or overfitted towards things that happen over and over again. This process improves the dataset's consistency by keeping only unique records. This makes it a more solid base for training forecast models.

b. Data Normalization: This is a standard step in preparing data that makes the features in a dataset fit into a common range of sizes without changing the way the different value ranges look. This is very important for datasets with different sensor readings, since devices may work and record data on different scales. Normalisation makes sure that every feature adds the same amount to the analysis. This way, no single feature can take over the model because of its size. It speeds up the training process and makes predictive models more accurate when this fair input is used.

c. Standard Scalar: To normalise data using the Standard Scalar method, the features are usually scaled so that their mean is zero and their standard deviation is one. When working with data that needs to be brought up to a standard scale so that machine learning algorithms can work best, this method is very helpful. The Standard Scalar method helps make features more consistent, which means that algorithms are less likely to favour features that are shown on bigger sizes. It also makes the program more sensitive to small differences in data, which is very important in situations like predicting machine failure, where accuracy is very important. The EDA shown in Figure 2 gives a picture of how the different sensor data points are spread out, which is very important for understanding the features of the dataset used to predict machine failure. When we start with "footfall," the data is much skewed, with most of the numbers being close to zero and only a few going up to higher counts. This means that most machine activities happen when there aren't many people around, but sometimes there are a lot of people. The 'tempMode' histogram has a bimodal distribution, which means that there are two usual operating temperature modes. These modes could mean different working conditions or shifts.

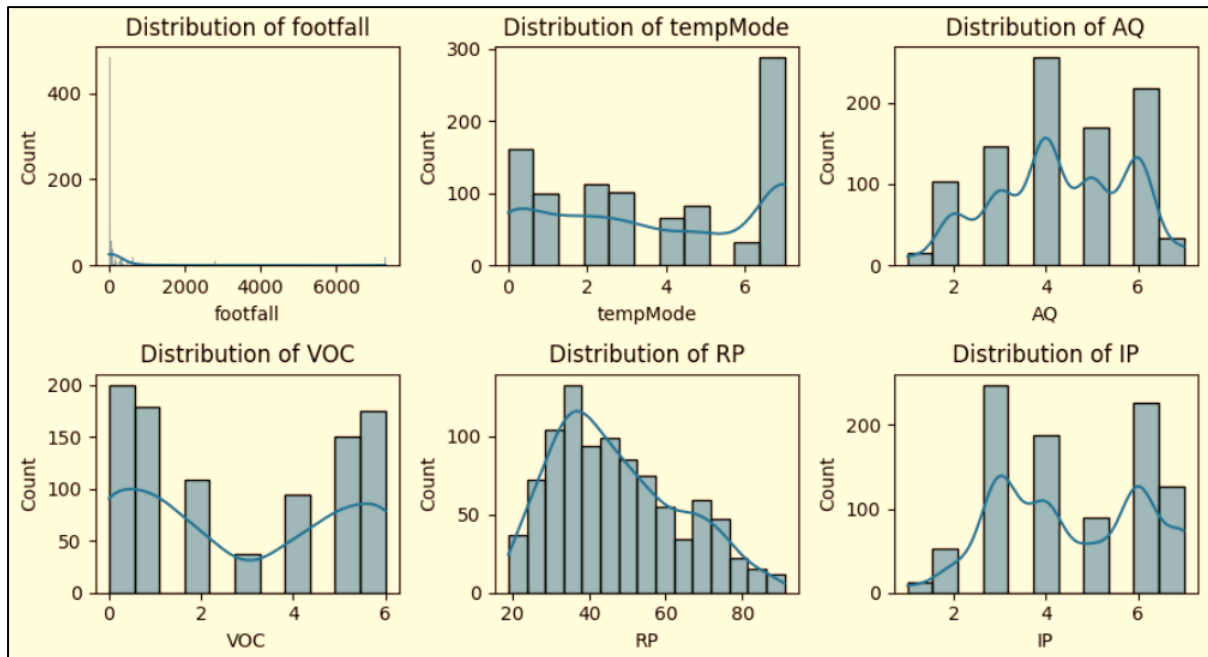


Figure 2. EDA of Dataset Features

The 'AQ' (Air Quality) data shows a fairly even spread of mid-range values, which suggests that air quality control is pretty consistent with some changes. 'VOC' (Volatile Organic Compounds) levels and 'IP' (Impact Pressure) levels both have mixed ranges, which mean that these devices may be used in a variety of situations. This could mean different steps of work or different settings for the machine. The 'RP' (Rotations Per Minute) graph displays a right-skewed distribution, meaning that lower rotation speeds are more common. However, there are enough cases of higher speeds to be noticed, which may be linked to higher rates of wear or failure. During the preparation step of data analysis, these distributions are very important because they help find outliers, figure out the data's core trends, and plan the next steps of normalization and feature engineering that will get the data ready for predictive modeling.

C. Feature Selection

a. Find the Correlation of Parameters with Label (FAIL): In this step, you count the correlation coefficients between each sensor parameter and the dataset's failure label. Correlation analysis helps figure out how different variables are related, showing which factors have the biggest effect on how likely it is that a machine will break down. It's important to know how well each trait can predict the future because that helps you build the model by focusing on the most important factors.

b. Select Feature with Correlation > 0.3: Once the correlations have been found, features with a failure name and a correlation value greater than 0.3 are chosen to train the model. This limit makes sure that only the traits that have a strong connection to the result are used, which makes the model more accurate and efficient. This chosen method makes the model simpler by leaving out data points that aren't important or have little effect, which could improve performance.

	fail
AQ	0.582895
USS	0.466712
VOC	0.797182
fail	1.000000

Figure 3. Selected Features

D. Data Balancing

Data balancing with the Synthetic Minority Over-sampling Technique (SMOTE) method is a very important part of predictive maintenance using the given dataset. When it comes to failure forecast, datasets are usually not fair, with a lot more instances of things that didn't fail than instances of things that did. When there is an unbalance like this, machine learning models may become biased towards the majority class, which makes them bad at predicting the minority (failure) class. This problem is fixed by SMOTE, which creates fake samples for the minority class based on the real examples. Interpolating between several similar minority class examples is what this method does. In this way, SMOTE helps to balance the dataset, which helps the forecast model, learn more about failures in general. This makes the model more accurate and reliable overall. Using SMOTE in the planning step makes sure that each class is represented equally during the training phase. This is important for getting results that are fair and correct. SMOTE Pseudo Algorithm for data Balancing:

Algorithm: SMOTE

Input:

S = Set of minority class samples

k = Number of nearest neighbors

N = Number of synthetic samples to generate per minority class sample

Output:

S' = Augmented set of minority class samples

Begin:

For each sample x_i in S do:

 Calculate distances to all other samples in S

 Identify k-nearest neighbors, denoted $NN_k(x_i)$

 For j from 1 to N do:

Randomly select $x_{\{in\}}$ from $NN_k(x_i)$

 Generate a random scalar λ between 0 and 1

 Compute new synthetic sample $x_{\{new\}}$:

$$x_{\{new\}} = x_i + \lambda * (x_{\{in\}} - x_i)$$

 Add $x_{\{new\}}$ to S'

Return S'
End

4. Machine Learning Classifier

a. Gradient Boosting Classifier:

The Gradient Boosting Classifier is a powerful ensemble machine learning algorithm that builds models in steps, like other boosting methods, and generalizes them by allowing optimization of an arbitrary differentiable loss function. By using more than one model in the final forecast, each new model gradually lowers the mistake of the older models. In gradient boosting, decision trees are trained one at a time, with each new tree helping to fix mistakes made by trees that were trained before it. It keeps getting better until no more changes can be made. Because it fixes the mistakes of the weaker models in the chain, it works especially well with datasets that aren't fair. Many people like this classifier because it's good at classifying things and doesn't overfit, even when the datasets are complicated and the relationships between the features aren't straight.

Gradient Boosting Classifier

Step 1: Initialization

- Initialize the model with a simple estimator that minimizes the loss function L:

$$F_{0(x)} = \arg \min_{\gamma} \sum L(y_i, \gamma)$$

Step 2: Iterative Boosting

- For $m = 1$ to M (where M is the number of boosting stages):

a. Compute the negative gradient (also called pseudo-residuals) of the loss function for each observation:

$$r_{\{im\}} = - \left[\frac{\partial L(y_i, F_{\{m-1\}(x_i)})}{\partial F_{\{m-1\}(x_i)}} \right]$$

b. Fit a new weak learner (e.g., a decision tree) to these residuals, denoting the output of this learner as $h_m(x)$:

$$h_{m(x)} = \text{learner}(r_{\{im\}}, x_i)$$

- c. Determine the best coefficient γ_m that fits $h_m(x)$ to $r_{\{im\}}$ using line search:

$$\gamma_m = \arg \min_{\gamma} \sum L(y_i, F_{\{m-1\}(x_i)} + \gamma h_{m(x_i)})$$

Step 3: Update the Model

- Update the model by adding the new learner scaled by the learning rate η :

$$F_m(x) = F_{\{m-1\}(x)} + \eta \gamma_m h_m(x)$$

Step 4: Final Model

- The final model $F_M(x)$ is the sum of the initial model and all the iteratively added models:

$$F_M(x) = F_{0(x)} + \eta \sum \gamma_m h_m(x)$$

End.

b. Isolation Forest:

Isolation Forest is an autonomous learning algorithm for finding anomalies that works by separating them, instead of modeling regular points like most other algorithms do. To use the method, random pick a feature and then pick a split value between the feature's highest and lowest values. This "isolates" the data. Since abnormalities are few and different, they are easier to identify compared to standard places. This results in shorter paths in the tree structure for irregularities during iterative splitting, making the method extremely efficient for big datasets. Isolation Forest is naturally suited to deal with outliers and is widely used in fraud detection, health tracking, and fault detection where quick separation of anomalies is important.

c. SVM (Support Vector Machine):

Support Vector Machine (SVM) is a guided machine learning method widely used for classification and regression problems. Essentially, it works by mapping data points in space and finding the hyperplane that best separates the classes. The vectors (data points) that form the hyperplane are the support vectors, and the size of the hyperplane depends on the features present in the dataset. People like SVM because it works well in spaces with a lot of dimensions and can handle both linear and non-linear boundaries, based on the kernel that is used (linear, polynomial, radial basis function, etc.). This makes SVM highly useful in areas like biology, picture recognition, and text categorization where the clear margin of difference is important.

d. Optimized Isolation Forest (Random Selection):

A better version of Isolation Forest is Optimized Isolation Forest (OIF). It does this by using Randomized Search to find the best values for things like the number of trees in the forest and the number of samples needed to make each tree. Randomised Search chooses parameters at random, looking at a wide range of numbers and giving you a useful way to find the best parameters for the model. This optimisation can help find anomalies more quickly by fine-tuning the model to fit the dataset's unique features better. When the parameter space is big and the distribution of outliers is complicated and hard to pick out using standard parameters, OIF works really well.

e. Hybrid 1D CNN – OIF:

As for the Hybrid 1D CNN–OIF, it takes the best parts of both Convolutional Neural Networks (CNNs) and Optimized Isolation Forests and puts them together in one structure. By using convolutional layers, 1D CNNs can pull out features from sets of data, which makes them perfect for time-series analysis. With an Optimized Isolation Forest added to the mix, the hybrid model uses CNNs' feature extraction skills to change the data before finding anomalies. This combined method works especially well when the data has trends in time or space that regular models might miss. This mix makes it easier to understand and find outliers, which makes it a powerful tool for dealing with large datasets where trends might not be clear at first but are needed to make predictions.

Model: "sequential_1"		
Layer (type)	Output Shape	Param #
conv1d_2 (Conv1D)	(None, 8, 32)	96
conv1d_3 (Conv1D)	(None, 7, 64)	4,160
flatten_1 (Flatten)	(None, 448)	0
dense_2 (Dense)	(None, 64)	28,736
dropout_1 (Dropout)	(None, 64)	0
dense_3 (Dense)	(None, 1)	65

Total params: 33,057 (129.13 KB)
Trainable params: 33,057 (129.13 KB)
Non-trainable params: 0 (0.00 B)

Figure 4. Hybrid 1D CNN – OIF Model Summary

f. Hybrid CNN – LSTM

It is a complex machine learning structure that combines the best features of Convolutional Neural Networks (CNNs) and Long Short-Term Memory networks (LSTMs) to handle patterns where changes in both space and time are significant. Through their convolutional layers, CNNs are great at getting hierarchical spatial features from raw data. This makes them perfect for handling data where understanding space is important, like pictures or sensor arrays. LSTMs are a type of recurrent neural network that is intended to find trends in long strings of data.

Model: "sequential_2"		
Layer (type)	Output Shape	Param #
conv1d_4 (Conv1D)	(None, 8, 32)	96
conv1d_5 (Conv1D)	(None, 7, 64)	4,160
lstm (LSTM)	(None, 64)	33,024
dropout_2 (Dropout)	(None, 64)	0
dense_4 (Dense)	(None, 1)	65

Total params: 37,345 (145.88 KB)
Trainable params: 37,345 (145.88 KB)
Non-trainable params: 0 (0.00 B)

Figure 5. Hybrid CNN – LSTM Model Summary

They can also see how things depend on each other over time and over long periods of time. The first step in this mixed model is to send data through CNN layers. These layers learn spatial structures and summarise them into more general features, which successfully reduces the number of dimensions in the data. The features that have been handled are then sent to LSTM layers, which look at the temporal aspects by keeping information over time. This makes it useful for time-series prediction, sequential data analysis, or any job where knowing the timing of events is very important. It is possible for the model to understand not only "where" something happens (thanks to CNNs), but also "when" it happens in connection to other events (thanks to LSTMs). In situations like video processing, where it's important to know both the spatial details of each frame and the time connections between frames, the combined CNN-LSTM model is very useful. In predictive maintenance, where sensor data may have both spatial arrangements and temporal sequences (for example, a set of readings over time), this model can accurately predict when equipment will break down by finding complex patterns in the data that models that only look at spatial or temporal patterns might miss. Predictive analytics works better and is more reliable when it can use both spatial and temporal data at the same time.

Hybrid CNN-LSTM Model

Step 1: Input Layer

- Input $x \in \mathbb{R}^d$, where d is the dimensionality of the input data (e.g., sensor readings, images).

Step 2: Convolutional Layer (CNN)

- Apply convolution operation:

$$z^l = w^l * x + b^l$$

- Where w^l and b^l are the weights and bias of the l -th convolutional layer, and $*$ denotes the convolution operation.

Step 3: Activation Function

- Apply a non-linear activation function (e.g., ReLU):

$$a^l = \max(0, z^l)$$

Step 4: Pooling Layer

- Apply pooling operation (e.g., max pooling) to reduce dimensionality:

$$p^l = \max(a^l_i), \text{ where } i \text{ defines the pooling region.}$$

Step 5: Long Short-Term Memory Layer (LSTM)

- Apply LSTM operations for sequence processing:

$$\begin{aligned} i_t &= \sigma(W_i * [h_{t-1}, x_t] + b_i) \\ f_t &= \sigma(W_f * [h_{t-1}, x_t] + b_f) \\ o_t &= \sigma(W_o * [h_{t-1}, x_t] + b_o) \\ g_t &= \tanh(W_c * [h_{t-1}, x_t] + b_c) \\ c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ h_t &= o_t \odot \tanh(c_t) \end{aligned}$$

- Where σ is the sigmoid activation, \tanh is the hyperbolic tangent activation, \odot denotes element-wise multiplication, and W terms and b terms are the weights and biases for respective gates.

Step 6: Output Layer

- Apply a fully connected layer to map LSTM output to desired output size:

$$y = W_h * h_t + b_y$$

- Where W_h is the weight matrix for the output layer, h_t is the output from the last LSTM cell, and b_y is the bias.

End

5. Apply XAI Technique

a. SHAP

SHAP, which stands for "SHapley Additive reasons," is a sturdy XAI (Explainable synthetic talent) technique that enables us recognizes how device mastering models make selections. For study tasks the usage of device learning classifiers like Gradient Boosting Classifier, Isolation woodland, SVM, and Hybrid CNN-LSTM models to expect when a device will break down, SHAP may be very useful in making the fashions more clean and reliable.

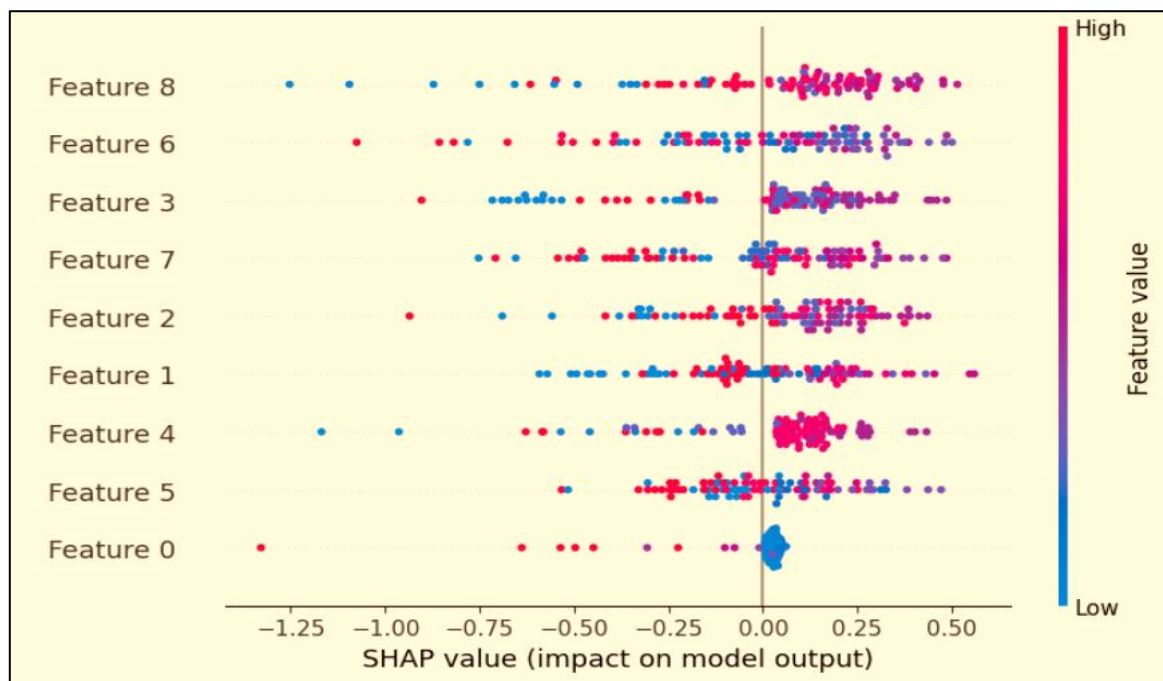


Figure 6. Feature Value

SHAP values come from the idea of game theory, specially Shapley values, which offer each trait quite a number that indicates how lots it enables with the estimate. This method breaks down a prediction into its elements to express how each factor affects it. This shall we researchers and engineers recognize how the readings from numerous gadgets have an effect

on the version's selection approximately whether or not a gadget will fail. To use SHAP in this case, you need to discern out the SHAP values for each trait across all the prediction fashions.

Algorithm: SHAP for XAI

Input:

Model = Trained Machine Learning Model (e.g., Gradient Boosting, SVM, CNN-LSTM)

Data = Feature Data used for Model Prediction

Output:

SHAP Values = Contributions of each feature to each prediction

Begin:

Step 1: Initialize the SHAP Explainer

- Create a SHAP explainer instance using the model.

Step 2: Calculate SHAP Values

- Compute SHAP values using the explainer for the entire dataset or a selected subset.

Step 3: Summarize the SHAP Values

- Calculate the mean absolute SHAP value for each feature across all data points to determine the average impact of each feature on the model output.

Step 4: Visualize the SHAP Values

- Generate visualization plots such as summary plots to display the distribution of the impacts each feature has on model predictions.

Step 5: Analyze High Impact Features

- Identify features with the highest mean absolute SHAP values as these are the most influential features in predicting the model's outcomes.

Step 6: Communicate Results

- Prepare reports or visualizations summarizing the SHAP values to communicate how each feature influences the model's predictions to stakeholders.

End

b. LIME

LIME, which stands for "Local Interpretable Model-agnostic Explanations," is an essential tool in the area of system learning. It makes complex predictive models used for predicting machine screw ups less complicated to recognize. In this take a look at, which uses advanced classifiers like Gradient Boosting, SVM, and complicated hybrid models like CNN-LSTM, LIME is very helpful because it indicates how those fashions determine what to do at the forecast degree for everyone. For every bet the model makes, LIME comes up with an answer through becoming it locally with a model that may be understood, like a choice tree or linear regression. To begin, the uncooked information is changed round a factor of hobby, and the version's estimates are checked to peer how they exchange. Then, these new samples are used to educate a model that is simpler to apprehend and is weighted by way of how close each altered point is to the authentic factor.

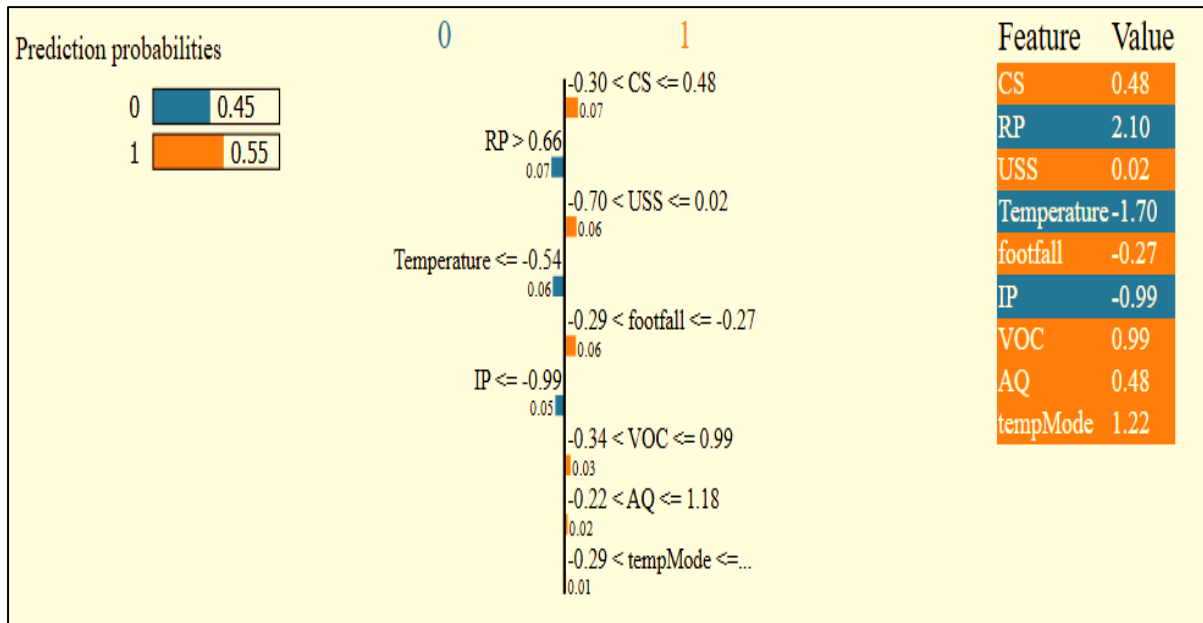


Figure 7. LIME Prediction probabilities

Researchers and practitioners can discover which developments are maximum important for making positive predictions via the usage of LIME. For instance, when looking to discern out why a machine will wreck down, LIME can show whether or not temperature, strain, or shaking are essential elements. This feature is mainly useful for first-rate-tuning the model to recognition greater precisely on key failure signs and symptoms. This improves the accuracy of predictions and builds trust within the model consequences among decision-makers and aid groups. This approach no longer solely takes the thriller out of complex fashions' "black containers," however it also helps with the method of improving and validating models time and again once more.

LIME Mathematical Model

1. Data Perturbation:

Let x be the original instance to be explained.

Generate a new dataset {z_1, z_2, ..., z_n} by perturbing x while keeping them close to x.

2. Proximity Measure:

Calculate the proximity $\pi_x(z)$ between the perturbed sample z and x using:

$$\pi_{x(z)} = \exp\left(-\frac{d(x, z)^2}{\sigma^2}\right)$$

where d(x, z) is the distance (e.g., Euclidean), and σ is the kernel width.

3. Model Fitting:

Use the dataset {z_1, z_2, ..., z_n} with model predictions {f(z_1), f(z_2), ..., f(z_n)}, and proximity weights { $\pi_x(z_1), \pi_x(z_2), \dots, \pi_x(z_n)$ } to fit a simple interpretable model g (like linear regression):

$$\text{Minimize over } g: \sum (\pi_{x(z_i)} * (f(z_i) - g(z_i))^2) + \Omega(g)$$

where $\Omega(g)$ is a complexity penalty to maintain the simplicity of g .

4. Explanation Generation:

The coefficients of model g indicate the importance of each feature in the prediction of x , providing an explanation for the model's behavior at x .

6. Result and Discussion

Figure 8 shows the training and validation loss and accuracy curves for a machine learning model over 50 epochs. It is most likely from one of the more advanced models used in studies on predicting machine failure, like the Hybrid CNN-LSTM. A big change can be seen in the Loss Curve. The training loss goes down sharply and then stays at a lower level. On the other hand, the validation loss goes down at first but then stops going down and starts fluctuating slightly. This could mean that the model is overfitting because it learns how to use what it has learnt from the training data but is having trouble applying it to new data it hasn't seen before.

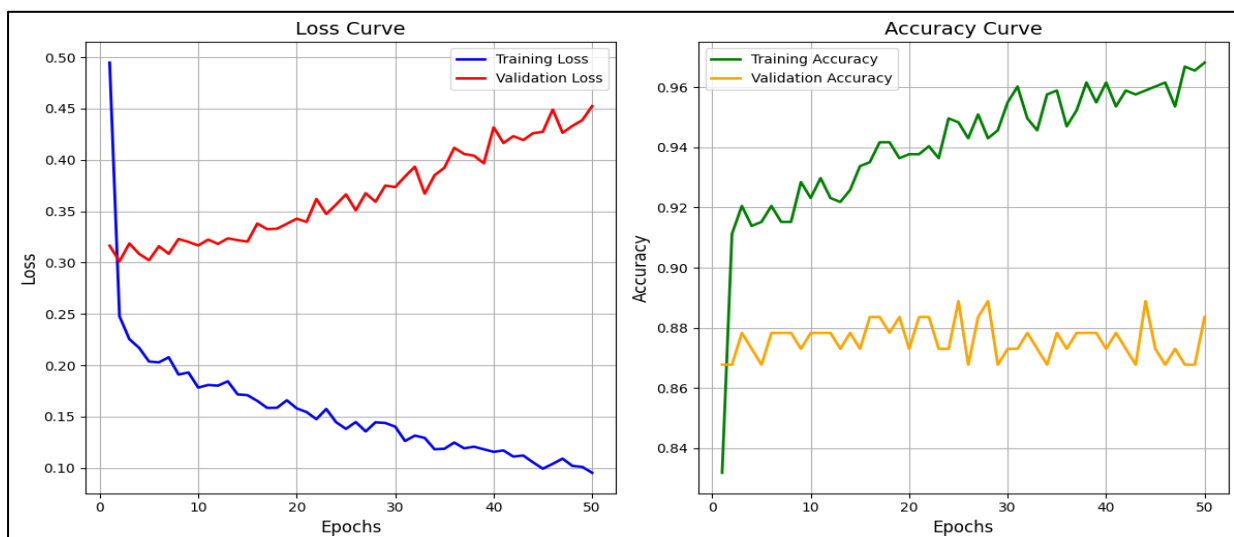


Figure 8. CNN – ISOMAP Accuracy and Loss Curve

The Accuracy Curve shows a big difference: training accuracy goes up constantly and reaches almost 96%, which means the model is getting better at using training data. On the other hand, the validation accuracy changes around 90% after an initial peak, showing that it's hard for the model to keep performing consistently on validation data. These trends show how important it is to keep an eye on both the training and validation metrics to make sure the model can be used in other situations and to avoid overfitting. These kinds of findings are very important for improving the reliability of predictive maintenance systems and fine-tuning model parameters.

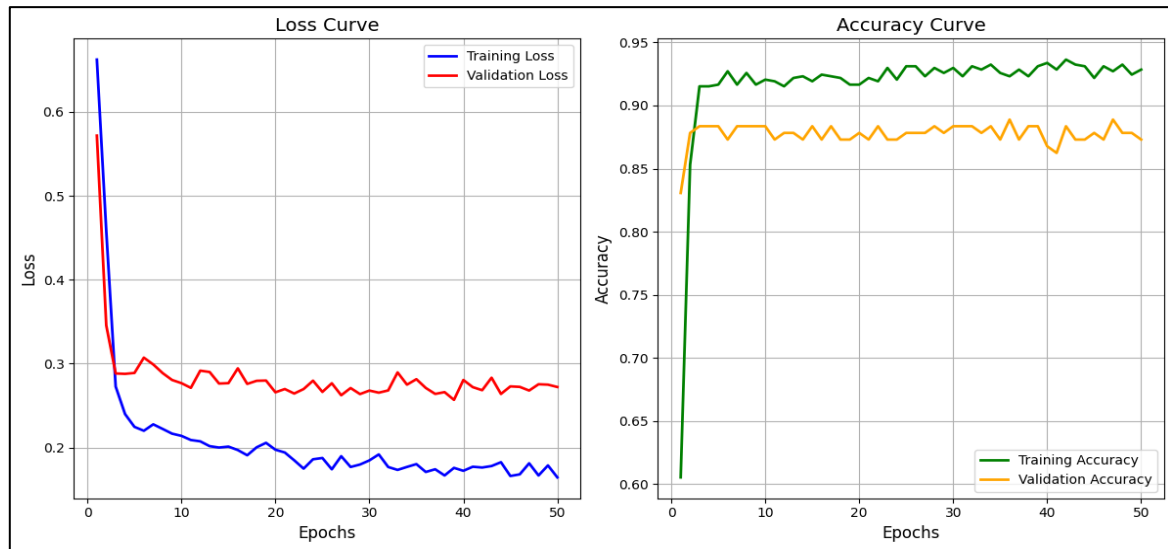


Figure 9. CNN – LSTM Accuracy and Loss Graph

Figure 9 shows how the training and evaluation process works for a machine learning model over 50 epochs. This model could be used to predict when a machine will break down. The loss curves show that the loss for both training and validation drops quickly and converges, which means that the model fits well without being too perfect. On the other hand, the accuracy graphs show that the training accuracy stays above 95%, while the validation accuracy stops around 85%. This means that the model works really well on training data but not so well on data it hasn't seen before. This shows where the model might need to be tweaked to make it more general.

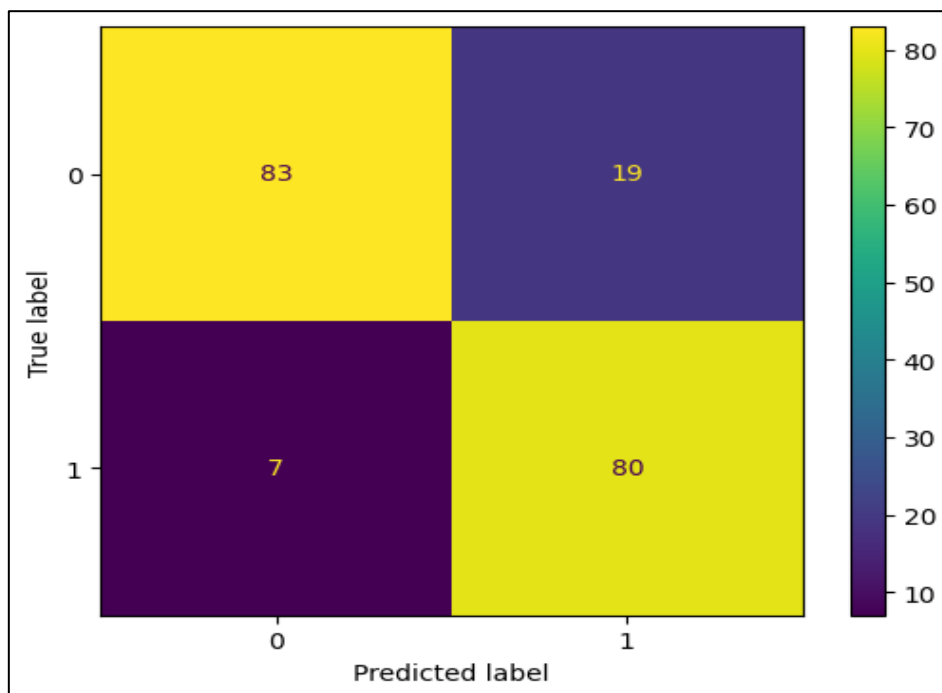


Figure 10. Confusion Matrix

The success of a classification model is shown by the confusion matrix. It shows that the model got '0' (no failure) 83 times and '1' (failure) 80 times right. 19 of the predictions were wrong, though; they said there would be a failure when there wasn't one, and 7 of the predictions were wrong and said there would be one. This means that the accuracy and sensitivity are pretty good, but there is still room for improvement in lowering the number of wrong guesses, as confusion matrix illustrates in figure 10. In Table 2, demonstrate comparison of how well different machine learning models did at predicting when machines would break down, with scores shown for accuracy and F1-score. The Gradient Boosting, Isolation Forest, and SVM models, which are some of the simplest, don't do very well. Gradient Boosting is a strong ensemble method that joins several weak prediction models to make a strong learner. It has a fair performance, with an F1-score of 52% for both accuracy and performance. This means that it does a good job of handling both the positive and negative classes, but it could be limited by how complicated or noisy the information is.

Table 2: Comparative Analysis

Model	Accuracy (%)	F1-Score (%)
Gradient Boosting	52	52
Isolation Forest	50	41
SVM	53	49
Optimized Isolation Forest (OIF)	86	87
Hybrid 1D CNN-OIF	87	87
Hybrid CNN-LSTM	92	88

The Isolation Forest model, with an F1-score of 41% and a 50% accuracy, performs the lowest of all. Its design is for locating anomalies rather than for balanced binary categorization, hence in projected maintenance scenarios, the failure class is rather crucial, so this significant F1-score decline makes it difficult to properly represent it. Known for performing well in high-dimensional settings, the SVM model only beats the Isolation Forest with 53% accuracy and 49% F1-score. Though SVM is a good predictor, the poor results might indicate issues like classes that cannot be separated linearly or insufficient space between classes in the feature space.

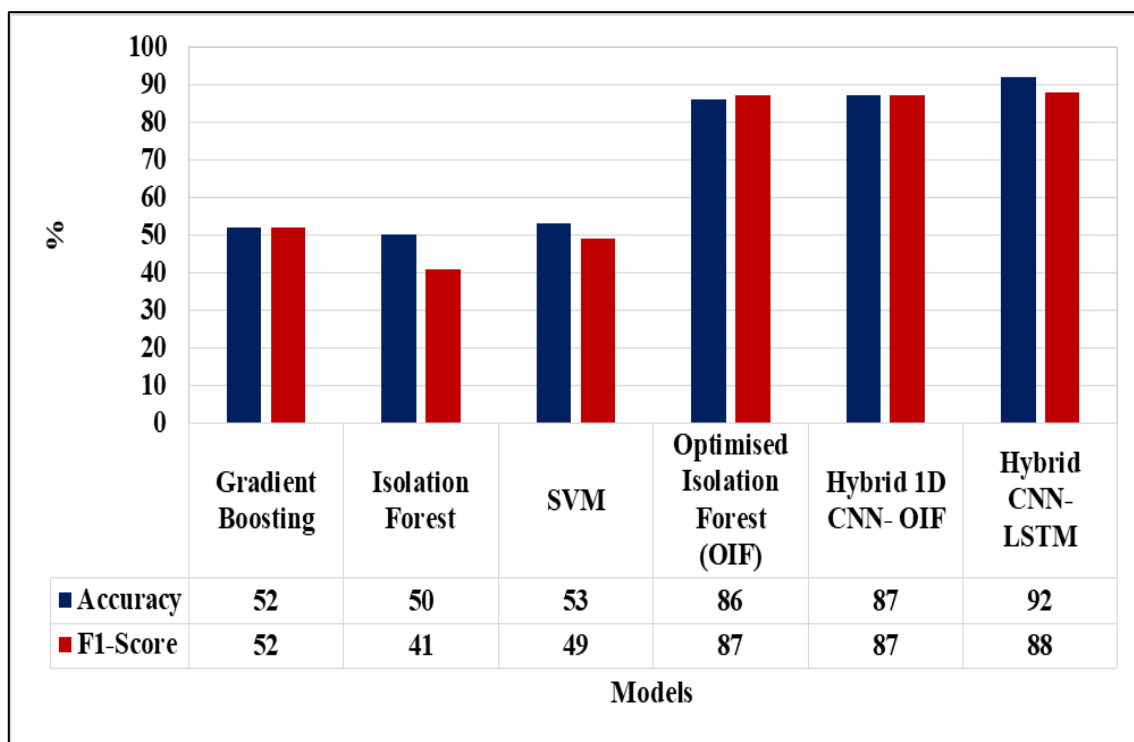


Figure 11. Comparative Analysis of Models

When we circulate on to more complicated models, the Optimized Isolation forest (OIF) is a massive breakthrough. It uses a Randomized search to fine-tune the settings of a regular Isolation forest. A score of 87% and an accuracy of 86% show that tuning and optimizing can significantly enhance the performance of a model, particularly in relation to tough responsibilities like predicting failure. The Hybrid 1D CNN-OIF model combines convolutional neural networks with Optimized Isolation woodland. It takes advantage of CNN's capacity to pull out beneficial features from time-series or linear facts, which is not unusual in sensor-pushed settings. An excellent 87% for both accuracy and F1-rating is accomplished by using this version, displaying how function extraction capabilities from CNNs and anomaly detection competencies from Isolation Forests paintings nicely together. The final version that stands out is the Hybrid CNN-LSTM model, which has an F1-score of 88% and an accuracy of 92%. This model combines the fantastic features of CNNs for extracting functions in space with LSTMs' strengths for processing records in a sequential order. This makes it best for complicated datasets that need to seize both spatial and temporal approaches for correct failure prediction. One clear trend that stands out from the table is that models get much better at predicting failures as they get more complicated and detailed, especially those that use mixed methods or optimized parameters, as illustrate in figure 11. In areas like predictive maintenance, where accuracy and quick spotting can avoid expensive downtimes and ensure operating efficiency, this means that spending money on advanced models and tuning strategies is very important.

Temp Mode:	0
AQ:	7
USS:	7
CS:	1
VOC:	36
Predict	
Predicted Fail Type : Fail	

Figure 12. Final Prediction Output

Figure 12 shows the user interface for an application called a prediction model that uses sensor information to guess when a machine will break down. Temperature Mode (Temp Mode), Air Quality (AQ), Ultrasonic Sensor Signal (USS), Current Sensor (CS), and Volatile Organic Compounds (VOC) are some of the factors that users can enter numbers for. After it put these numbers, the model figures out how likely it is that the machine will break down. Based on the numbers that were given, the model suggests a loss ("Fail"). This tool is useful for preventative maintenance because it lets workers step in and fix problems before they happen, which improves the stability and efficiency of operations.

7. Conclusion

The study showed a deep look into how different machine learning models can accurately guess what will happen and explain how they make decisions in predictive maintenance situations. Many tests and studies were conducted to evaluate how well models such as Gradient Boosting, Isolation Forest, SVM, Optimized Isolation Forest (OIF), Hybrid 1D CNN-OIF, and Hybrid CNN-LSTM performed. Our findings indicate that while simple models like Gradient Boosting, Isolation Forest, and SVM are beneficial for understanding predictive maintenance, more sophisticated models are more accurate and dependable than they are. Speed-wise, the Optimized Isolation Forest and mixed models using deep learning have shone out. With the greatest accuracy and dependability of all the models tested, the Hybrid CNN-LSTM model which combines the spatial feature identification capability of CNNs with the sequential data processing capacity of LSTMs was the most effective. This work has also improved by means of XAI techniques as SHAP and LIME, which provide us additional insight on how sophisticated artificial intelligence models operate without really seeing them. These techniques have not only clarified how certain characteristics influence model outcomes but also contributed to the clarity and dependability of these models. LIME's localized explanations, for instance, have proved very useful in determining which characteristics most influence the projections of the model. This has resulted in concentrated modifications and an improved knowledge of how things operate in the actual world. This study shows how AI is

changing in predictive maintenance, where advanced models combined with methods for explanation can greatly enhance the decision-making process. For predictive analytics to reach its full potential, it will be important for businesses to combine complex AI models with strong XAI systems as they continue to go digital. This not only makes sure that operations run smoothly and with little downtime, but it also helps people trust and understand automatic systems. In the future, both AI and XAI will need to keep getting better to keep up with the growing complexity of real-world data and decision-making situations.

References

- [1] Kashani, Z.A.; Pakzad, R.; Fakari, F.R.; Haghparast, M.S.; Abdi, F.; Kiani, Z.; Talebi, A.; Haghgoo, S.M. Electromagnetic fields exposure on fetal and childhood abnormalities: Systematic review and meta-analysis. *Open Med.* 2023, 18, 20230697.
- [2] Bosch-Capblanch, X.; Esu, E.; Dongus, S.; Oringanje, C.M.; Jalilian, H.; Eysers, J.; Oftedal, G.; Meremikwu, M.; Rössli, M. The effects of radiofrequency electromagnetic fields exposure on human self-reported symptoms: A protocol for a systematic review of human experimental studies. *Environ. Int.* 2022, 158, 106953.
- [3] Martin, S.; De Giudici, P.; Genier, J.C.; Cassagne, E.; Doré, J.F.; Ducimetière, P.; Evrard, A.S.; Letertre, T.; Ségala, C. Health disturbances and exposure to radiofrequency electromagnetic fields from mobile-phone base stations in French urban areas. *Environ. Res.* 2021, 193, 110583.
- [4] Ramirez-Vazquez, R.; Escobar, I.; Vandenbosch, G.A.; Arribas, E. Personal exposure to radiofrequency electromagnetic fields: A comparative analysis of international, national, and regional guidelines. *Environ. Res.* 2024, 246, 118124.
- [5] Aerts, S.; Deprez, K.; Verloock, L.; Olsen, R.G.; Martens, L.; Tran, P.; Joseph, W. RF-EMF Exposure near 5G NR Small Cells. *Sensors* 2023, 23, 3145.
- [6] Liu, S.; Tobita, K.; Onishi, T.; Taki, M.; Watanabe, S. Electromagnetic field exposure monitoring of commercial 28-GHz band 5G base stations in Tokyo, Japan. *Bioelectromagnetics* 2024, 45, 281–292.
- [7] Kovalnogov, V.N.; Fedorov, R.V.; Generalov, D.A.; Chukalin, A.V.; Katsikis, V.N.; Mourtas, S.D.; Simos, T.E. Portfolio insurance through error-correction neural networks. *Mathematics* 2022, 10, 3335.
- [8] Alonso Robisco, A.; Carbó Martínez, J.M. Measuring the model risk-adjusted performance of machine learning algorithms in credit default prediction. *Financ. Innov.* 2022, 8, 1–35.
- [9] Lundberg, S.M.; Erion, G.; Chen, H.; DeGrave, A.; Prutkin, J.M.; Nair, B.; Katz, R.; Himmelfarb, J.; Bansal, N.; Lee, S.I. From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* 2020, 2, 56–67.
- [10] Kshirsagar, R.; Hsu, L.Y.; Greenberg, C.H.; McClell, M.; Mohan, A.; Shende, W.; Tilmans, N.P.; Guo, M.; Chheda, A.; Trotter, M.; et al. Accurate and Interpretable Machine

- Learning for Transparent Pricing of Health Insurance Plans. Proc. AAAI Conf. Artif. Intell. 2021, 35, 15127–15136.
- [11] Duddalwar, A., Khobragade, P. (2025). An Optimization of Healthcare Operation Management Using Machine Learning. In: Bhateja, V., Chakravarthy, V.V.S.S.S., Anguera, J., Ghosh, A., Flores Fuentes, W. (eds) Signal Processing, Telecommunication and Embedded Systems with AI and ML Applications. ICMEET 2023. Lecture Notes in Electrical Engineering, vol 1281. Springer, Singapore. https://doi.org/10.1007/978-981-97-8422-6_36
- [12] Du, Y.; Rafferty, A.R.; McAuliffe, F.M.; Wei, L.; Mooney, C. An explainable machine learning-based clinical decision support system for prediction of gestational diabetes mellitus. *Sci. Rep.* 2022, 12, 1170.
- [13] Islam, M.R.; Ahmed, M.U.; Barua, S.; Begum, S. A systematic review of explainable artificial intelligence in terms of different application domains and tasks. *Appl. Sci.* 2022, 12, 1353.
- [14] Clement, T.; Kemmerzell, N.; Abdelaal, M.; Amberg, M. XAIR: A Systematic Metareview of Explainable AI (XAI) Aligned to the Software Development Process. *Mach. Learn. Knowl. Extr.* 2023, 5, 78–108.
- [15] Zhang, S.; Meng, X.; Bhagia, S.; Ji, A.; Dean Smith, M.; Wang, Y.; Liu, B.; Yoo, C.G.; Harper, D.P.; Ragauskas, A.J. 3D Printed Lignin/Polymer Composite with Enhanced Mechanical and Anti-Thermal-Aging Performance. *Chem. Eng. J.* 2024, 481, 148449.
- [16] Mosi, G.; Ikua, B.W.; Kabini, S.K.; Mwangi, J.W. Characterization and Modeling of Mechanical Properties of Additively Manufactured Coconut Fiber-Reinforced Polypropylene Composites. *Adv. Mater. Phys. Chem.* 2024, 14, 95–112.
- [17] Anerao, P.; Kulkarni, A.; Munde, Y. A Review on Exploration of the Mechanical Characteristics of 3D-Printed Biocomposites Fabricated by Fused Deposition Modelling (FDM). *Rapid Prototyp. J.* 2023, 30, 430–440.
- [18] Anerao, P.; Kulkarni, A.; Munde, Y.; Shinde, A.; Das, O. Biochar Reinforced PLA Composite for Fused Deposition Modelling (FDM): A Parametric Study on Mechanical Performance. *Compos. Part C Open Access* 2023, 12, 100406.
- [19] Rendas, P.; Figueiredo, L.; Cláudio, R.; Vidal, C.; Soares, B. Investigating the Effects of Printing Temperatures and Deposition on the Compressive Properties and Density of 3D Printed Polyetheretherketone. *Prog. Addit. Manuf.* 2023, 1–17.
- [20] Ahmad, M.N.; Ishak, M.R.; Mohammad Taha, M.; Mustapha, F.; Leman, Z. Irianto Mechanical, Thermal and Physical Characteristics of Oil Palm (*Elaeis guineensis*) Fiber Reinforced Thermoplastic Composites for FDM—Type 3D Printer. *Polym. Test.* 2023, 120, 107972.
- [21] Fisher, T.; Almeida, J.H.S., Jr.; Falzon, B.G.; Kazancı, Z. Tension and Compression Properties of 3D-Printed Composites: Print Orientation and Strain Rate Effects. *Polymers* 2023, 15, 1708.