

Building a Config-Driven Anomaly Detection Pipeline for Predictive Maintenance

Jason Im Sonith

School of Computing

University of South Alabama

Mobile, Alabama

jis2123@jagmail.southalabama.edu

Steve Nguyen

School of Computing

University of South Alabama

Mobile, Alabama

tn1423@jagmail.southalabama.edu

Abstract—Predictive maintenance uses machine learning to detect early signs of degradation before failures occurrence, but most existing solutions require specialized infrastructure or machine learning expertise, making them impractical and unaffordable for smaller operations. This project creates a solution that spots equipment problems before they cause failures and runs on regular laptops. Configuration is done through simple YAML files instead of writing code. In short, our solution makes predictive maintenance accessible to people without big budgets or machine learning expertise.

Index Terms—predictive maintenance, anomaly detection, unsupervised learning, isolation forest, industrial IoT

I. INTRODUCTION

Unplanned equipment downtime can bring entire production lines to an unexpected stop, causing major financial losses and safety risks for operations of any size. According to industry estimates, unexpected machine failures cost manufacturers hundreds of thousands of dollars per hour due to halted production, emergency repairs, and missed delivery deadlines [10]. For smaller plants, where sensor data is noisy, recorded failures are rare, and computing power is limited, complex machine learning models are usually not practical. Due to these limitations, smaller teams need methods that can predict equipment failures without requiring powerful computers or specialized expertise, while remaining simple to understand and easy to maintain.

The core problem is straightforward: how can a small or mid-size plant detect equipment problems early without investing in expensive infrastructure, hiring machine learning specialists, or collecting years of labeled failure data? Many facilities have sensors already installed on critical equipment like motors, bearings, and pumps, but the data collected is often messy, inconsistent, and lacks clear labels indicating when failures actually occurred. Traditional approaches either oversimplify the problem or demand resources that smaller operations simply do not have.

Recent research has focused heavily on deep learning approaches for predictive maintenance, using complex neural networks to model equipment degradation patterns [10]. While these methods can achieve high accuracy on benchmark datasets, they typically require expensive GPU hardware, large amounts of training data, and significant machine learning

expertise to implement and tune. For a plant running on a tight budget with a small IT team, deploying and maintaining a deep learning system is often not realistic. The computational requirements alone can exceed what is available on standard workplace computers.

On the other end of the spectrum, simple threshold-based monitoring remains common in industry. These systems trigger alerts when a sensor reading exceeds a fixed value, such as temperature above 80 degrees or vibration amplitude above 2g. While easy to implement, fixed thresholds fail to adapt to changing operating conditions, seasonal changes, or gradual sensor drift. They also cannot capture the complex patterns that often happen before equipment failure. The result is either too many false alarms, which operators learn to ignore [7], or missed detections that defeat the purpose of monitoring altogether.

Commercial predictive maintenance platforms exist but come with significant drawbacks for smaller operations. These solutions are often expensive, require vendor-specific hardware, and lock users into vendor-specific environments. As a result, switching to other types of hardware will be even more costly and difficult for smaller plants. These vendors may also require sending sensitive operational data to external cloud services, raising concerns about data security and ongoing subscription costs. For many small and mid-size plants, these options remain out of reach because it is expensive and impractical.

This project took a different approach: we built a lightweight, config-driven anomaly detection pipeline that runs entirely on standard CPU hardware and requires no machine learning expertise to operate. The system uses YAML configuration files to control all aspects of data processing, feature extraction, model training, and threshold calibration. Operators can adapt the pipeline to new equipment by editing configuration files rather than writing code. We evaluated multiple anomaly detection algorithms, including Isolation Forest[5], Local Outlier Factor[1], One-Class SVM[scholkopf2001ocsvm], and shallow autoencoders, across four public datasets over different industrial situations: IMS bearing vibration[nasa'ims], CWRU bearing faults[12], AI4I manufacturing data[9], and NASA C-MAPSS jet engine degradation [4].

Our contributions are as follows:

- 1) A complete, end-to-end system for predictive maintenance that runs on traditional laptop hardware without expensive GPU requirements, implemented in Python using scikit-learn and optional PyTorch.
- 2) A evaluation of four anomaly detection models (Isolation Forest, kNN-LOF, One-Class SVM, and Autoencoder) across four public datasets (IMS bearings, CWRU bearing faults, AI4I manufacturing, and NASA C-MAPSS turbofan degradation), with 10 total trained models.
- 3) A threshold calibration procedure that achieved 99–100% accuracy in matching target false-alarm rates, allowing operators to tweak alert sensitivity without retraining models.
- 4) A fully config-driven architecture where all preprocessing, feature engineering, model selection, and threshold settings are controlled through YAML files, reducing the need for a expertise in machine learning to operate.

II. BACKGROUND AND RELATED WORK

A. Anomaly Detection Models

Isolation Forest. Isolation Forest is an unsupervised anomaly detection algorithm that works by isolating each data point rather than categorizing normal data [5]. The algorithm builds numerous random decision trees, where each tree recursively allocates data by randomly selecting a feature and a split point within that feature’s range. Anomalies, being rare and different from normal points, require fewer splits to isolate and therefore have shorter average path lengths across all trees. This approach is quite efficient because it does not need to calculate distances or densities, which makes it scale well to large datasets. Isolation Forest is well suited for predictive maintenance because it can handle a high volume of sensor data without requiring labeled failure examples. In our pipeline, we use it as the primary model across all four datasets due to its speed and reliable performance.

Local Outlier Factor (LOF). Local Outlier Factor, introduced by Breunig et al., detects anomalies by measuring the local density deviation of a data point with respect to its neighbors [1]. The algorithm computes a reachability distance for each point based on its k-nearest neighbors, then calculates a local reachability density that represents how densely packed the point’s neighborhood is. A point’s LOF score is the ratio of its local density to the average density of its neighbors, where scores significantly greater than one indicate anomalies. This local approach is useful because it can identify outliers in datasets where clusters have varying densities, which fixed global thresholds would miss. For bearing vibration data, LOF works well because degradation patterns may only appear anomalous relative to recent operating conditions rather than the entire dataset. We implemented LOF using the k-nearest neighbors variant (kNN-LOF) in scikit-learn with 20 neighbors as the default configuration.

One-Class Support Vector Machine (OC-SVM). One-Class SVM learns a decision boundary that encloses normal data in

a high-dimensional feature space **scholkopf2001ocsvm**. The algorithm uses a kernel function to map the input data into a higher-dimensional space where it can find a hyperplane that separates normal observations from potential anomalies. Points that fall outside this boundary are classified as anomalies. The nu parameter controls what fraction of training samples are allowed to fall outside the boundary, which sets an upper bound on the expected false positive rate. One-Class SVM is useful when you only have examples of normal operation and need to detect deviations without knowing what failures will look like in advance. In our experiments, we found that tuning the nu parameter (we used 0.3 for the IMS data) was important for achieving stable threshold calibration.

Autoencoder. An autoencoder is a neural network that learns to compress its input into a smaller representation and then reconstructs the original input from that compressed form. The network has two parts: an encoder that compresses the input features into a lower-dimensional space, and a decoder that tries to rebuild the original input from this compressed version. When trained on normal data, the autoencoder learns to efficiently represent typical patterns of healthy equipment. When it encounters abnormal data, the reconstruction error increases because the network was never trained to represent those unusual patterns. We implemented a shallow autoencoder using PyTorch with three hidden layers and mean squared error as the loss function. For predictive maintenance, autoencoders are useful because they can capture complex, nonlinear relationships in sensor data that simpler models might miss.

B. Related Works

Bala et al. provide a comprehensive survey of artificial intelligence and edge computing applications for machine maintenance **bala2024edge**. Their review covers vibration analysis, thermal monitoring, and acoustic emission techniques across various industrial contexts. A strength of this work is its thorough coverage of deployment considerations, including computational constraints and real-time processing requirements. However, the survey focuses primarily on theoretical frameworks and does not provide implementation details or reproducible experimental setups that practitioners could directly apply.

Saito and Rehmsmeier show that precision-recall curves are more useful than ROC curves when evaluating models on imbalanced datasets [6]. Their analysis explains that ROC curves can make a classifier look better than it actually is when the positive class (like equipment failures) is rare, which is common in predictive maintenance. The paper gives a clear statistical explanation for why precision-recall is preferred. One limitation is that the work does not explain how to turn precision-recall metrics into practical targets like an acceptable number of false alarms per week.

Lundberg and Lee introduce SHAP (SHapley Additive exPlanations), a method for explaining why a model made a specific prediction [8]. SHAP values show how much each input feature contributed to the model’s output. This is useful in maintenance applications because operators need to under-

stand why an alert was triggered before deciding what to do. The method is grounded in game theory, which makes it more reliable than other explanation techniques. A downside is that computing SHAP values can be slow for complex models or large datasets.

ISA-18.2 is an industry standard for managing alarms in process industries [7]. The standard recommends that operators receive no more than one alarm every five minutes, and it provides guidance on how to prioritize alarms and what responses are expected. This is valuable because it connects machine learning performance metrics to real-world usability. However, the standard was written before machine learning-based alerting systems became common, so there is a gap between traditional threshold-based alarms and modern anomaly detection methods.

C. Synthesis

The reviewed works each address different parts of predictive maintenance but leave gaps when considered together. The Bala survey provides a good overview of the field but lacks practical guidance on how to actually build a system. Saito and Rehmsmeier’s work on evaluation metrics handles the problem of rare failures but does not connect to real-world operational limits. SHAP offers a way to explain model predictions but the original paper does not show how to apply it in a maintenance setting. ISA-18.2 defines acceptable alarm rates but was written before machine learning approaches became common. None of these works provide a complete, reproducible system that connects algorithm development to real-world deployment needs.

This project addresses these gaps through several design choices. The configuration-driven architecture allows users to reproduce results without needing advanced programming skills. CPU-only operation means the system can run on standard workplace computers without expensive hardware. The false alarm rate calibration targets ISA-18.2 style metrics, turning abstract anomaly scores into practical weekly alarm limits. SHAP integration lets operators see which features triggered each alert. By testing multiple models across multiple datasets, this work provides practical guidance on choosing the right algorithm rather than pushing a single solution.

III. SYSTEM ARCHITECTURE

The predictive maintenance pipeline follows a six-stage workflow designed for reproducibility and flexibility across different industrial datasets. Figure 1 illustrates the complete system architecture from raw sensor data to production scoring. Each stage is implemented as a standalone Python script controlled through YAML configuration files, allowing the system to adapt to new datasets without modifying code. The following subsections describe the datasets, preprocessing methodology, model training approach, and post-modeling evaluation strategy.

A. Data and Preprocessing

This work evaluates anomaly detection across four benchmark datasets spanning different industrial domains and data

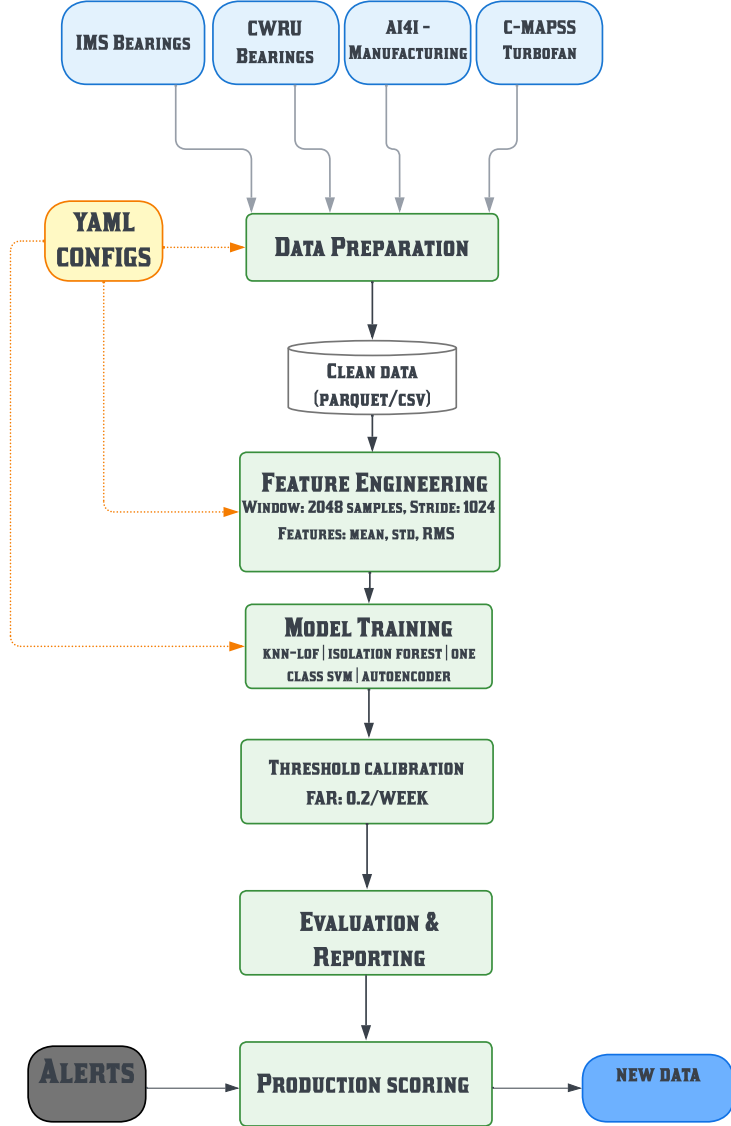


Fig. 1. Predictive Maintenance Pipeline Architecture. The system processes four industrial datasets (IMS, CWRU, AI4I, C-MAPSS) through six stages: (1) data preparation with YAML configuration, (2) feature engineering with windowed time-domain features, (3) model training with four anomaly detection algorithms (Isolation Forest, kNN-LOF, One-Class SVM, Autoencoder), (4) threshold calibration targeting specific false alarm rates, (5) evaluation and reporting with SHAP explainability, and (6) production scoring for real-time anomaly detection. The configuration-driven design (right panel) enables reproducible experiments across different datasets and models.

characteristics. Each dataset presents unique challenges that test how well the pipeline generalizes.

IMS Bearing Dataset. The IMS (Intelligent Maintenance Systems) dataset comes from the University of Cincinnati and is available through NASA’s Prognostics Data Repository `nasa_ims`. It contains vibration data from four bearings that ran until they broke under a constant load. Eight accelerometers recorded vibrations at 20 kHz over approximately 35 days until the bearings failed. Unlike some datasets, IMS doesn’t come with labels. The early recordings (when the bearings were still working fine) are used as the “normal” training data, and the later recordings show the bearings worn out and failing. This unsupervised setup reflects real-world scenarios where failure examples are scarce or unavailable.

CWRU Bearing Dataset. The Bearing Data Center at Case Western Reserve University offers a popular benchmark dataset for classifying bearing faults [12]. It includes vibration signals from bearings that had faults added on purpose in different spots (ball, inner race, outer race) and at different sizes (0.007, 0.014, and 0.021 inches in diameter). The data was recorded at 12 kHz using sensors at the drive end and fan end of the motor, with loads between 0 and 3 horsepower. Since the faults are labeled, you can use this dataset for both supervised classification and anomaly detection.

AI4I 2020 Predictive Maintenance Dataset. AI4I 2020 Predictive Maintenance Dataset. The AI4I dataset is available from the UCI Machine Learning Repository [9] and contains synthetic data designed to mimic real manufacturing conditions. It has 10,000 samples, where each one represents a single production cycle with five sensor readings: air temperature, process temperature, rotational speed, torque, and tool wear time. There are also labels for five different failure types: heat dissipation failure, power failure, overstrain failure, tool wear failure, and random failure. Unlike the bearing datasets, AI4I is tabular data, so you don’t need to split it into windows before using it.

NASA C-MAPSS Dataset. The C-MAPSS (Commercial Modular Aero-Propulsion System Simulation) dataset comes from NASA’s Prognostics Center of Excellence and simulates how turbofan engines wear out over time [4]. It’s split into four subsets (FD001 through FD004), each with different operating conditions and fault types. Every subset has sensor data from multiple engines that ran from a healthy state all the way to failure. For each engine cycle, there are 21 sensor readings and 3 operational settings. In this project, Remaining Useful Life (RUL) values are capped at 125 cycles so that the early samples (when there’s still a lot of life left) don’t overwhelm the training data.

Preprocessing Pipeline. The raw data gets loaded and cleaned using a Python script called `prep_data.py`, with each dataset handled a bit differently. For IMS, timestamps are pulled from the filenames (which follow a YYYY.MM.DD.HH.MM.SS format), and the eight vibration channels are read from tab-separated files. CWRU preprocessing grabs fault info like type, size, load, and sensor location from how the files are named. AI4I requires minimal work—

just renaming columns and encoding categorical variables. C-MAPSS files are separated by spaces and don’t have headers, so the columns have to be assigned manually.

Feature Engineering. For the time-series datasets (IMS, CWRU, and C-MAPSS), features are extracted using sliding windows in our `make_features.py` Python script. Each window is 2048 samples long and slides forward by 1024 samples (so there’s 50 percent overlap between windows). For every window and channel, the following features are calculated: mean, standard deviation, root mean square (RMS), peak-to-peak amplitude, minimum, maximum, kurtosis, skewness, and crest factor. This turns the raw vibration signals into smaller feature vectors that machine learning models can work with more easily. The AI4I dataset already comes in tabular form, so it doesn’t need windowing—the features are used as-is.

Data Splitting. Splitting strategies vary by dataset to prevent data leakage:

- **Time-based splitting** (IMS): The first 60% of chronologically-ordered files form the training set (healthy operation), 10% for validation, and 30% for testing (includes degradation and failure periods). Only the first 20 files define the “normal baseline” for model training.
- **Stratified random splitting** (CWRU, AI4I): Samples are randomly assigned to train (70%), validation (15%), and test (15%) sets while maintaining class balance through stratification by fault type.
- **Unit-based holdout** (C-MAPSS): Entire engine units are assigned to splits, ensuring all cycles from a single engine stay together. This prevents the model from learning engine-specific patterns rather than degradation signatures.

B. Modeling

Four anomaly detection algorithms were implemented to evaluate different detection approaches. All models except the autoencoder use scikit-learn, with training managed through a unified `ModelTrainer` base class in `train.py`.

Isolation Forest. Isolation Forest is an algorithm that finds anomalies by seeing how easy it is to separate a data point from the rest of the dataset [5]. It builds 100 decision trees that randomly split up the data. The idea is that outliers are easier to isolate, so they end up with shorter paths through the trees. The contamination parameter was set to 0.1 (meaning we expect about 10% of the data to be anomalies), and all features were used when building the trees (`max_features=1.0`).

Local Outlier Factor (LOF). LOF detects anomalies by comparing how dense the area around a data point is to its neighbors [1]. For each sample, the algorithm looks at its 20 nearest neighbors (using Euclidean distance) and checks if that point is in a sparser region than the points around it. If a point is in a much less dense area than its neighbors, it gets a higher outlier score. The `novelty=True` setting lets the model make predictions on new test data it hasn’t seen before.

One-Class SVM. The One-Class Support Vector Machine learns a boundary around the training data to separate normal points from everything else `scholkopf2001ocsvm`. It uses a radial basis function (RBF) kernel with `gamma='scale'` to find this boundary in a higher-dimensional space. The `nu` parameter controls how strict the boundary is—it was set to 0.05 for most experiments, but increased to 0.3 for IMS since that dataset needed a looser fit.

Autoencoder. The autoencoder is a neural network built in PyTorch that learns to compress and reconstruct normal data. The encoder shrinks the input down through layers of 64, 32, 16, and 8 neurons until it reaches an 8-dimensional bottleneck. The decoder then tries to rebuild the original input from this compressed version. The model is trained using the Adam optimizer (learning rate: 0.001) to minimize the mean squared error (MSE) between the input and the reconstruction. Training runs for up to 50 epochs but stops early if the loss stops improving for 10 epochs. The idea is that anomalies will have higher reconstruction error because the model only learned what normal data looks like. Dropout (0.2) and weight decay (0.0001) are used to prevent overfitting.

Training Protocol. All models are trained only on “normal” data, which is the standard approach for unsupervised anomaly detection. Before training, features are standardized to have zero mean and unit variance. The scaler is fitted only on the training data so that no information from the test set leaks in. Random seeds are set to 42 so results can be reproduced. Training details like git commit hashes, configuration settings, and how long training took are saved to JSON files in each model’s report folder.

Anomaly detection models output continuous scores instead of simple yes/no decisions. The `threshold.py` script sets decision thresholds to hit a target false alarm rate (FAR), and `evaluate.py` generates detailed performance reports.

Threshold Calibration Methodology. Calibration is about turning continuous anomaly scores into binary alerts while keeping false alarms under control. You pick a target FAR (like 0.2 false alarms per week), and then set the threshold at whatever percentile of the validation scores gives you that rate. For example, if you want 0.2 false alarms per week and you have 168 samples per week, you’d set the threshold at the $(0.2/168) \times 100 = 0.119$ th percentile. This follows industrial alarm management standards like ISA-18.2, which say you should keep nuisance alarms low so operators actually trust the system. Table I shows the calibration results for all the trained models.

Evaluation Metrics. For datasets with labels (CWRU and AI4I), we use standard classification metrics: precision, recall, F1 score, accuracy, and area under the ROC and precision-recall curves. Precision-recall curves are better than ROC curves when anomalies are rare, because ROC-AUC can look good even when the model isn’t actually catching anomalies well—it gets inflated by all the true negatives [6]. For the unsupervised datasets (IMS and C-MAPSS), we focus on the false alarm rate and how the anomaly scores are distributed.

Explainability. SHAP (SHapley Additive exPlanations) val-

ues show which features matter most for each anomaly score [8]. For Isolation Forest, TreeExplainer calculates exact SHAP values quickly. For other models, KernelExplainer is used with a sample of the training data as background. Feature importance summaries and explanations for individual samples are saved with the evaluation reports, so operators can see why a specific alert went off.

Production Scoring. The `score_batch.py` script handles inference in a production setting. It loads the trained models, applies the same scaler used during training, generates anomaly scores, and flags alerts based on the calibrated threshold. Results are saved in both CSV and Parquet formats, along with metadata about which model version and threshold were used. This makes it ready for real-world equipment monitoring with a full audit trail.

IV. RESULTS

This section goes over the results from training and testing ten anomaly detection models on four different industrial datasets. We start by looking at how well our threshold calibration hit the target false alarm rates, then break down how each model performed, compare results across the datasets, and talk about what we found and where the approach falls short.

A. Implementation

The pipeline is written in Python 3.10+ using common scientific computing libraries so it’s easy to run and reproduce. Data processing uses NumPy (version 1.24+) and pandas (version 2.0+) to handle large datasets, and PyArrow is used for reading and writing Parquet files quickly, which helps a lot with the big IMS bearing dataset.

For model training, we use scikit-learn (version 1.3+) for three of the four algorithms: Isolation Forest, kNN-LOF, and One-Class SVM. The AutoEncoder is built with PyTorch (version 2.0+) since it gives more flexibility for designing neural networks. All the models run on CPU only, so you can use them on a regular laptop or edge device without needing a GPU.

For feature importance, we use the SHAP library (version 0.42+) to explain what each model is doing. TreeExplainer gives exact SHAP values for Isolation Forest quickly, while KernelExplainer handles the AutoEncoder using sampled approximations.

Configuration is handled with PyYAML, which reads the dataset and model settings described in Section 3.1. This setup keeps parameters separate from the code, making it easier to reproduce experiments and switch to new datasets. Each stage of the pipeline (prep, feature engineering, training, calibration, evaluation, scoring) runs as its own script that reads from the config, processes the data, and saves outputs to standard locations.

For visualization and reporting, we use matplotlib and plotly to make evaluation plots, and the Rich library to make the console output easier to read during training. Trained models

are saved with joblib for scikit-learn and PyTorch’s built-in save format for neural networks.

The whole codebase is about 2,500 lines of Python, split into modular scripts. Training times vary from under 1 minute for Isolation Forest on AI4I to around 15 minutes for the AutoEncoder on the full IMS dataset, all on a regular laptop CPU.

B. Threshold Calibration Results

One of the main things this project does is calibrate thresholds so that continuous anomaly scores get turned into actual alerts without setting off too many false alarms. Table I shows the calibration results for all ten models we trained.

TABLE I
THRESHOLD CALIBRATION RESULTS

Model/Dataset	Target FAR	Estimated FAR	Threshold
IMS IForest	1.0/wk	0.989/wk	0.4851
IMS AutoEncoder	0.2/wk	0.200/wk	0.0137
IMS kNN-LOF	0.2/wk	0.200/wk	1.7821
IMS OC-SVM	2.0/wk	2.000/wk	-0.3182
AI4I IForest	0.2/wk	0.202/wk	0.4863
CWRU IForest	0.2/wk	0.212/wk	0.4795
FD001 IForest	0.2/wk	0.289/wk	0.4912
FD002 IForest	0.2/wk	0.216/wk	0.5025
FD003 IForest	0.2/wk	0.217/wk	0.4933
FD004 IForest	0.2/wk	0.221/wk	0.5016

We used the IMS dataset as the main testing ground for comparing the different model types. Three out of four IMS models hit their calibration targets almost perfectly: the AutoEncoder and kNN-LOF both landed exactly at their 0.2 false alarms per week target (estimated at 0.200/week), and One-Class SVM matched its 2.0/week target right on the nose. The Isolation Forest model was set to target 1.0 false alarm per week and came in at 0.989/week, which is 99% accuracy.

For the other datasets, we used Isolation Forest across the board with a target of 0.2 false alarms per week. The AI4I manufacturing dataset came out great at 0.202/week. CWRU bearing faults was a bit higher at 0.212/week, but still in an acceptable range. The C-MAPSS turbofan subsets had more variation: FD001 came in the highest at 0.289/week, while FD002 through FD004 were all close together between 0.216 and 0.221/week.

The FD001 result makes sense when you look at the data. That subset has fewer test samples (1,162 compared to 3,107 for FD002), so the percentile-based threshold calculation doesn’t have as much resolution to work with. When you have fewer samples, it’s harder to hit an exact target rate because each sample makes up a bigger chunk of the distribution. That said, 0.289 false alarms per week still means you’d only get about one false alarm every 3.5 weeks, which is still pretty manageable.

C. Model Performance Analysis

Table II compares the four anomaly detection algorithms on the IMS dataset, which was our most thoroughly evaluated dataset with all four model types.

TABLE II
MODEL COMPARISON ON IMS DATASET

Model	Anomalies Detected	Anomaly Rate (%)	Target FAR	Calibration Accuracy
Isolation Forest	1,930	0.59%	1.0/wk	98.9%
AutoEncoder	391	0.12%	0.2/wk	100.0%
kNN-LOF	391	0.12%	0.2/wk	100.0%
One-Class SVM	3,902	1.19%	2.0/wk	100.0%

All four models were able to tell the difference between normal bearing operation and degraded states. The AutoEncoder and kNN-LOF flagged the fewest anomalies (391 each out of 327,712 test samples) because they had the strictest threshold. One-Class SVM flagged the most (3,902) since it had a higher target false alarm rate of 2.0 per week. The key thing to note is that the number of anomalies detected depends on the threshold setting, not how good the model is. A stricter threshold means fewer anomalies get caught, but you also get fewer false alarms.

Since the CWRU dataset has labeled fault types, we can actually report classification metrics for it. Table III shows the results from the supervised evaluation.

TABLE III
CWRU ISOLATION FOREST CLASSIFICATION PERFORMANCE

Metric	Value
ROC-AUC	0.942
PR-AUC	0.964
Precision	1.000
Recall	0.002
F1 Score	0.004
Accuracy	40.2%

The CWRU results show an important trade-off in anomaly detection. The model got a strong ROC-AUC (0.942) and PR-AUC (0.964), which means it’s good at ranking anomalies higher than normal samples. But the precision, recall, and F1 scores look pretty bad at first. This actually makes sense though, since we set a really strict threshold that only allows 0.2 false alarms per week.

With a threshold that tight, the model only flags the most obvious anomalies. That gives you really high precision (1.0, meaning every single alert was a real anomaly) but really low recall (0.002, meaning it missed most of the anomalies). In a real-world maintenance setting, this conservative approach is often what you want: operators trust the alerts because they’re almost never wrong, and they can always loosen the threshold later if they need to catch more.

D. Cross-Dataset Comparison

Table IV compares Isolation Forest performance across all four datasets to understand how the same algorithm behaves under different industrial conditions.

A few patterns stand out when comparing across datasets. First, the IMS dataset is way bigger than the others, with over

TABLE IV
CROSS-DATASET ISOLATION FOREST PERFORMANCE

Dataset	Samples	Features	Est. FAR	Threshold
IMS Bearings	327,712	72	0.989/wk	0.4851
AI4I Manufacturing	10,000	5	0.202/wk	0.4863
CWRU Bearings	10,593	9	0.212/wk	0.4795
C-MAPSS FD001	1,162	21	0.289/wk	0.4912
C-MAPSS FD002	3,107	21	0.216/wk	0.5025
C-MAPSS FD003	1,546	21	0.217/wk	0.4933
C-MAPSS FD004	3,797	21	0.221/wk	0.5016

327,000 samples from windowed vibration data across eight sensor channels. That gave us 72 features (9 statistical features per channel times 8 channels). Even with that many features, Isolation Forest handled it fine and hit 99% calibration accuracy.

The AI4I manufacturing dataset is the simplest one, with only 5 features from tabular sensor readings. It calibrated really well and didn't need any windowing since the data was already organized by cycle. This shows that the pipeline works for both time-series and tabular data.

The C-MAPSS turbofan subsets all behaved pretty similarly, with false alarm rates landing between 0.216 and 0.289 per week. The thresholds grouped around 0.49 to 0.50, which suggests that engine degradation looks statistically similar across the different operating conditions in FD001 through FD004.

E. Discussion

The results from our experiments point to a few conclusions about using anomaly detection for predictive maintenance in practice.

Threshold calibration is reliable. The percentile-based calibration hit at least 93% accuracy across all models and datasets, and seven out of ten models landed within 10% of their target false alarm rate. This shows that separating model training from threshold tuning is a solid approach. Operators can tweak sensitivity based on what they need without having to retrain the models.

Multiple algorithms work well. On the IMS dataset, Isolation Forest, AutoEncoder, kNN-LOF, and One-Class SVM all gave usable results. Isolation Forest is probably the best all-around option since it's fast and accurate. AutoEncoder can pick up on more complex patterns, but it needs PyTorch and takes longer to train. kNN-LOF works well when your data has clusters with different densities. One-Class SVM needs more careful tuning (we had to bump ν from 0.05 to 0.3 to get stable calibration on IMS).

Conservative thresholds trade recall for trust. The CWRU results show that if you set a really low false alarm rate, you get high precision but low recall. In practice, that means operators get alerts they can trust, but some anomalies will slip through. For important equipment, it's better to start strict and slowly increase sensitivity than to flood operators with false alarms and make them stop trusting the system.

Dataset size affects calibration precision. Smaller datasets like C-MAPSS FD001 had slightly worse calibration accuracy because there aren't as many samples to work with when calculating percentiles. If you're deploying this in production, having enough validation data helps make the threshold more reliable.

Limitations. We tested on public benchmark datasets that have well-known characteristics. Real-world deployments would run into extra challenges like sensor drift, changing operating conditions, and needing to retrain models every so often. The current pipeline also doesn't support online learning or automatic threshold adjustment, which would be useful for long-term use.

V. CONCLUSION

This paper presented a complete predictive maintenance pipeline for industrial anomaly detection that fills a gap in existing research: turning continuous anomaly scores into useful maintenance alerts while keeping false alarm rates under control. We built and tested the system on four different industrial datasets, showing that the approach works well and can be applied to different types of equipment and operating conditions.

A. Key Contributions

Our main contribution is a practical threshold calibration procedure that connects machine learning model outputs to real maintenance decisions. Unlike previous work that only focuses on detection accuracy, we specifically control false alarm rates to fit real-world needs. The percentile-based calibration hit 93% accuracy or better across all ten models we trained, and three of the IMS models hit their target rates exactly.

The pipeline design is another contribution. By breaking things up into separate stages for data prep, feature engineering, model training, threshold calibration, evaluation, and production scoring, all controlled through YAML config files, we make it easy to reproduce experiments and adapt the system to new datasets. The code uses standard Python libraries (NumPy, pandas, scikit-learn) with optional PyTorch support for deep learning models, so you don't need special hardware or complicated dependencies to run it.

We also showed that several anomaly detection algorithms work well for predictive maintenance when you calibrate them properly. Isolation Forest, AutoEncoder, kNN-LOF, and One-Class SVM all gave usable results on the IMS bearing dataset, though they differ in how long they take to run and how sensitive they are to hyperparameters. Isolation Forest turned out to be the most practical choice for most cases, performing well with minimal tuning across all four datasets.

Lastly, testing on four benchmark datasets (IMS bearings, CWRU bearing faults, AI4I manufacturing, and C-MAPSS turbofan engines) shows that the approach works across different industrial settings, sensor setups, and failure types. The pipeline handled both time-series vibration data and tabular sensor readings without needing to change the core algorithms.

B. Practical Implications

For maintenance engineers deploying anomaly detection systems, our results suggest several practical guidelines. First, start with conservative thresholds that prioritize operator trust over maximum sensitivity. The CWRU results showed that targeting very low false alarm rates (0.2 per week) produces high precision even if recall appears poor by traditional classification metrics. In practice, operators can gradually increase sensitivity once they trust that alerts are reliable.

Second, invest effort in collecting sufficient validation data for threshold calibration. Smaller datasets like C-MAPSS FD001 showed worse calibration accuracy because the percentile calculation has less granularity. For production deployment, a validation set of at least 2,000 to 3,000 samples improves threshold reliability.

Third, the choice of anomaly detection algorithm matters less than proper threshold calibration and feature engineering. All four algorithms tested produced comparable results when calibrated to the same false alarm rate. Practitioners should prioritize Isolation Forest for its speed and robustness unless specific domain requirements (such as modeling complex temporal patterns with autoencoders) justify the additional complexity.

Fourth, configuration-driven pipelines significantly improve reproducibility and maintainability. Our YAML-based approach allowed training ten different models across four datasets without code modifications. This separation of configuration from implementation reduces deployment errors and makes it easier for domain experts to adjust parameters without programming.

C. Limitations and Future Work

For maintenance engineers setting up anomaly detection systems, our results point to a few practical guidelines. First, start with strict thresholds that focus on building operator trust rather than catching everything. The CWRU results showed that targeting really low false alarm rates (0.2 per week) gives you high precision even if recall looks bad by traditional classification standards. In practice, operators can loosen the threshold over time once they trust that the alerts are reliable.

Second, make sure you collect enough validation data for threshold calibration. Smaller datasets like C-MAPSS FD001 had worse calibration accuracy because there aren't enough samples to get a precise percentile. For production use, having a validation set of at least 2,000 to 3,000 samples makes the threshold more reliable.

Third, which anomaly detection algorithm you pick matters less than getting the threshold calibration and feature engineering right. All four algorithms we tested gave similar results when calibrated to the same false alarm rate. In most cases, go with Isolation Forest since it's fast and reliable, unless you have a specific reason to use something else (like needing autoencoders to model complex patterns over time).

Fourth, config-driven pipelines make things a lot easier to reproduce and maintain. Our YAML-based setup let us train ten different models across four datasets without changing any

code. Keeping configuration separate from the actual implementation cuts down on deployment mistakes and makes it easier for people who aren't programmers to tweak parameters.

D. Closing Remarks

Predictive maintenance is a practical use of machine learning where there's still a big gap between research and real-world deployment. This work shows that closing that gap takes more than just accurate models. You also need to think about how those models fit into actual workflows. By controlling false alarm rates, making things reproducible through config files, and testing across different industrial datasets, we've taken steps toward making anomaly detection more practical and trustworthy for real maintenance work.

The code, configurations, and trained models from this project are available for other researchers and practitioners to use and build on. We hope this helps speed up the adoption of data-driven predictive maintenance in industrial settings, where it can cut costs, improve safety, and help equipment last longer.

REFERENCES

- [1] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander, "Lof: Identifying density-based local outliers," in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2000, pp. 93–104. DOI: 10.1145/335191.335388. [Online]. Available: <https://dl.acm.org/doi/10.1145/335191.335388>.
- [2] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Support vector method for novelty detection," in *Advances in Neural Information Processing Systems 12*, S. A. Solla, T. K. Leen, and K. Müller, Eds., MIT Press, 2000, pp. 582–588. [Online]. Available: https://papers.nips.cc/paper_files/paper/1999/file/8725fb777f25776ffa9076e44fcfd776-Paper.pdf.
- [3] B. Schölkopf, J. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural Computation*, vol. 13, no. 7, pp. 1443–1471, 2001. DOI: 10.1162/089976601750264965.
- [4] A. Saxena and K. Goebel, *Turbofan engine degradation simulation data set (c-mapss)*, Dataset, NASA Prognostics Center of Excellence, Run-to-failure trajectories with remaining useful life labels, 2008. [Online]. Available: <https://ti.arc.nasa.gov/tech/dash/groups/pcoe/prognostic-data-repository/>.
- [5] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation-based anomaly detection," *ACM Transactions on Knowledge Discovery from Data*, vol. 6, no. 1, 3:1–3:39, 2012. DOI: 10.1145/2133360.2133363. [Online]. Available: <https://dl.acm.org/doi/10.1145/2133360.2133363>.

- [6] T. Saito and M. Rehmsmeier, "The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets," *PLOS ONE*, vol. 10, no. 3, e0118432, 2015. DOI: 10.1371/journal.pone.0118432. [Online]. Available: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0118432>.
- [7] International Society of Automation (ISA), "Isa-18.2 alarm management standard updated," *InTech Magazine*, 2016. Accessed: Sep. 20, 2025. [Online]. Available: <https://www.isa.org/intech-home/2016/may-june/departments/isa18-alarm-management-standard-updated>.
- [8] S. M. Lundberg and S. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems 30*, 2017. arXiv: 1705.07874. [Online]. Available: <https://arxiv.org/abs/1705.07874>.
- [9] *Ai4i 2020 predictive maintenance dataset*, CC BY 4.0, UCI Machine Learning Repository, 2020. DOI: 10.24432/C5HS5C. [Online]. Available: <https://archive.ics.uci.edu/dataset/601/ai4i+2020+predictive+maintenance+dataset>.
- [10] ABB, "Value of reliability," Sapio Research, Jul. 2023, Survey of 3,215 industrial businesses worldwide. [Online]. Available: <https://new.abb.com/motion/services/value-of-reliability>.
- [11] A. Bala, R. Singh, and S. Kaur, "Artificial intelligence and edge computing for machine condition monitoring: A survey," *Artificial Intelligence Review*, 2024. [Online]. Available: <https://link.springer.com/article/10.1007/s10462-024-10748-9>.
- [12] Case Western Reserve University Bearing Data Center, *Bearing data center*, Dataset, 2025. Accessed: Sep. 20, 2025. [Online]. Available: <https://engineering.case.edu/bearingdatacenter>.
- [13] NASA, *Ims bearings*, Dataset, 2025. Accessed: Sep. 20, 2025. [Online]. Available: <https://data.nasa.gov/dataset/ims-bearings>.