

## **INTEGRATING MATHEMATICAL MODELING AND DEEP LEARNING FOR UNCERTAINTY-AWARE FAULT DIAGNOSIS IN INDUSTRIAL ROTATING MACHINERY**

**Primawati<sup>1</sup>, Ferra Yanuar<sup>2\*</sup>, Dodi Devianto<sup>3</sup>, Remon Lapisa<sup>4</sup>, Fazrol Rozi<sup>5</sup>, Arda Yunianta<sup>6</sup>**

Department of Mechanical Engineering, Faculty of Engineering, Universitas Negeri Padang, Indonesia<sup>14</sup>

Department of Mathematics and Data Science, Universitas Andalas, Indonesia<sup>1235</sup>

Department of Information Technology, Politeknik Negeri Padang, Indonesia<sup>5</sup>

Department of Information Systems, Faculty of Computing and Information Technology, King Abdulaziz University, Rabigh, Saudi Arabia<sup>6</sup>

primawati@ft.unp.ac.id, ferrayanuar@sci.unand.ac.id\*, ddevianto@sci.unand.ac.id,

remonlapisa@ft.unp.ac.id, fazrol@pnp.ac.id, ayunianta@kau.edu.sa

Received: 10 October 2025, Revised: 12 March 2026, Accepted: 04 April 2026

\*Corresponding Author

---

### **ABSTRACT**

*In Industry 4.0, reliable fault diagnosis is critical for minimizing downtime and preventing catastrophic failures in rotating machinery. However, conventional deep learning models often operate deterministically, lacking the ability to quantify prediction uncertainty—a limitation that hinders risk-based maintenance decisions. This study aims to develop a hybrid deep learning framework that integrates Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), and Bayesian inference for uncertainty-aware fault diagnosis. The model extracts spatial features from Short-Time Fourier Transform (STFT) spectrograms via CNN, models temporal dynamics from raw vibration signals via LSTM, and quantifies prediction uncertainty using Monte Carlo Dropout ( $T=50$ ). Evaluated on the benchmark Case Western Reserve University (CWRU) bearing dataset with an 80/20 data partitioning under six operating conditions, the hybrid architecture achieves an accuracy of 99.14% and an F1-score of 0.9914, significantly outperforming standalone CNN (97.42%) and LSTM (84.12%) models. The integration of probabilistic inference enhances decision reliability by providing confidence estimates for each prediction. This work contributes a robust, uncertainty-aware model that effectively captures both spatial and temporal patterns, offering significant implications for safety-critical industrial predictive maintenance systems.*

**Keywords:** Fault Diagnosis, Rotating Machinery, Hybrid CNN-LSTM, Monte Carlo Dropout, Uncertainty Quantification

### **1. Introduction**

In the era of Industry 4.0, predictive maintenance has become a critical strategy for ensuring the reliability and efficiency of rotating machinery in manufacturing, energy, and transportation systems (Wahid et al., 2022; Syah et al., 2026). Unplanned equipment failures not only disrupt operations but also incur significant downtime and repair costs (Seoni et al., 2023; Yuan et al., 2020). Consequently, condition-based monitoring through vibration analysis has emerged as a key enabler of early fault detection (Jia & Sharma, 2021; Mohd Ghazali & Rahiman, 2021; Primawati et al., 2025; Tama et al., 2023; Rozi et al., 2025). However, as industrial systems become more complex, the demand for diagnostic models has shifted from mere accuracy to reliability and interpretability (Raj et al., 2024; Song et al., 2019). In safety-critical environments, false alarms can lead to unnecessary maintenance costs, while missed detections can cause catastrophic failures (Primawati, Yanuar, et al., 2026; Tama et al., 2023; Thoppil & Vasu, 2025). Therefore, uncertainty quantification is essential to enable risk-based maintenance decisions and dynamic threshold adaptation under variable load conditions. (Abdar et al., 2021).

Traditional diagnostic methods rely on signal processing techniques such as Fast Fourier Transform (FFT), Root Mean Square (RMS), and kurtosis to identify anomalies in time-domain or frequency-domain representations (Ventricci et al., 2024; Wang et al., 2022). While these approaches are well-established, they often fail to capture complex, non-stationary patterns

present in real-world vibration signals, particularly during incipient fault stages (Choi & Joe, 2024; Zhan et al., 2022). To overcome these limitations, machine learning (ML) and deep learning (DL) models have gained prominence. Support Vector Machines (SVMs) and ensemble methods like Gradient Boosting have demonstrated over 90% accuracy in failure prediction tasks (Rafati & Shaker, 2024; Rihi et al., 2024). However, their performance depends heavily on manual feature engineering, which limits generalization across diverse operational conditions (Shahin et al., 2023).

Recent advances in deep learning have significantly advanced vibration-based fault diagnosis by enabling automatic feature extraction from raw sensor data (Abdar et al., 2021; Stroescu & Olcay, 2022; Zhan et al., 2022; Lokeshwaran et al., 2025). Convolutional Neural Networks (CNNs) are highly effective in extracting spatial patterns from vibration signals, particularly when signals are transformed into time–frequency representations such as spectrograms (Chen et al., 2020; Lubis et al., 2025; Pinedo-Sánchez et al., 2020). To leverage both capabilities, hybrid CNN-LSTM architectures have gained prominence in the last five years. For instance, (Wahid et al., 2022) introduced an attention-based CNN-LSTM model for fault localization, and (Borré et al., 2023) proposed a hybrid framework for machine failure prediction. Despite these advances, a critical research gap remains: most existing hybrid frameworks operate deterministically. They provide a class label without indicating the confidence level of the prediction. This overconfidence makes such models sensitive to noise and domain shifts, lacking the interpretability required to distinguish between genuine faults and transient anomalies (Yanuar et al., 2023; Desnelita et al., 2025; Purbojo & Wijaya, 2025).

This lack of uncertainty quantification hinders the adoption of AI in risk-based maintenance decisions. In industrial settings, knowing *when* the model is uncertain is as critical as the prediction itself (Thoppil & Vasu, 2025). Uncertainty quantification is essential for: (1) minimizing false alarm costs by flagging low-confidence predictions for human review, (2) enabling dynamic threshold adaptation under variable load conditions, and (3) supporting risk-aware maintenance scheduling. Bayesian neural networks, particularly those using Monte Carlo (MC) Dropout, address this by estimating prediction uncertainty (Huang et al., 2024). Following the seminal work of (Gal & Ghahramani, 2016), who interpreted dropout as approximate Bayesian inference, applying MC Dropout during inference allows the model to quantify predictive variance. This probabilistic framework provides not only class predictions but also confidence estimates, which is vital for reliability-aware systems.

Furthermore, while Bayesian neural networks using Monte Carlo (MC) Dropout have been proposed in other domains, there is a distinct lack of existing probabilistic hybrid CNN–LSTM models specifically applied to vibration-based fault diagnosis (Ali et al., 2020). Existing studies often utilize benchmark datasets like CWRU without acknowledging challenges such as load variability and fault severity imbalance, or demonstrating how their models cope with the resulting uncertainty. Unlike existing deterministic CNN-LSTM fault classifiers, this study integrates Bayesian inference via Monte Carlo Dropout to quantify predictive uncertainty, thereby enhancing reliability in safety-critical industrial applications (Fang et al., 2020, 2020; Kendall & Gal, 2017; Fauzan et al., 2025).

To validate this approach, the study utilizes the Case Western Reserve University (CWRU) bearing dataset. While widely recognized as a benchmark, the CWRU dataset presents specific theoretical challenges, including load variability and fault severity imbalance, which test the robustness of diagnostic algorithms against domain shifts. Therefore, this study addresses these gaps by proposing a Hybrid Bayesian CNN-LSTM model. The research is guided by the following explicit objectives and questions:

- a. Performance Superiority: Does the Bayesian CNN–LSTM significantly outperform deterministic hybrids in terms of accuracy and generalization under variable loads?
- b. Uncertainty Utility: Does uncertainty estimation improve fault confidence and reduce risk in decision-making processes, particularly under noisy conditions?
- c. Feature Synergy: How effectively does the fusion of spatial (STFT) and temporal (raw signal) features distinguish structurally similar faults?

This work contributes to the state-of-the-art in three analytical ways:

- a. Methodological Novelty: It presents the first integration of Bayesian inference into a

CNN–LSTM architecture specifically for bearing fault diagnosis, moving beyond deterministic predictions to probabilistic reliability.

- b. Reliability-Aware Maintenance: It demonstrates how uncertainty quantification can support risk-aware decision-making, potentially reducing unnecessary maintenance interventions and false alarms.
- c. Industrial Deployment Perspective: The model is designed with computational efficiency in mind (e.g., MC Dropout  $T=50$ ), facilitating potential embedding into edge devices for real-time monitoring systems. This addresses key Industry 4.0 requirements regarding latency, sensor fusion, and scalability in resource-constrained environments.

The model is rigorously evaluated using the CWRU bearing dataset under multiple load levels, including normal, inner race, ball, outer race, and high-load conditions. Results show that the hybrid architecture achieves exceptional performance (99.14% accuracy) while maintaining low validation loss (0.0867), demonstrating superior generalization and robustness compared to standalone models. By bridging the gap between high-accuracy deep learning and probabilistic safety standards, this research offers a mathematically grounded, uncertainty-aware solution for modern predictive maintenance.

## 2. Literature Review

### 2.1 Scope and Selection Methodology

To ensure a comprehensive and current theoretical foundation, a systematic literature review was conducted focusing on publications from 2020 to 2025. Searches were performed across major databases (IEEE Xplore, ScienceDirect, Springer) using keywords such as "*fault diagnosis*," "*rotating machinery*," "*hybrid deep learning*," "*uncertainty quantification*," and "*predictive maintenance*." The selection prioritized studies that addressed industrial AI systems, reliability engineering, and real-time feasibility. This review critically synthesizes recent advances in signal processing, deep learning architectures, and probabilistic methods, identifying key limitations that motivate the proposed Bayesian CNN–LSTM framework.

### 2.2 Traditional and Machine Learning Approaches

Early fault diagnosis relied heavily on signal processing techniques such as Fast Fourier Transform (FFT) and statistical indicators (e.g., RMS, kurtosis) to detect periodic impulses associated with bearing defects (Li et al., 2022; Wang et al., 2022). While effective for stationary signals, these methods struggle with non-stationary industrial environments where load variations obscure fault signatures (Cariño et al., 2020). To improve robustness, machine learning (ML) classifiers such as Support Vector Machines (SVM) and Random Forests were applied to handcrafted features (Primawati et al., 2025). Although ensemble methods like Stacked SMO achieved accuracies up to 98.3%, they remain limited by manual feature engineering, which requires extensive domain expertise and reduces scalability across different machinery types (Stroescu & Olcay, 2022). Furthermore, traditional ML models often lack the capacity to generalize under domain shift, such as changes in motor load or speed, leading to performance degradation in real-world deployments (Rafati & Shaker, 2024).

Recent studies have addressed class imbalance challenges in industrial monitoring systems through advanced sampling and loss function strategies, demonstrating improved performance in quality monitoring applications (Asrol & Pratama, 2025; Thoppil & Vasu, 2025; Wahid et al., 2022). However, these approaches still require labeled fault data, which is often scarce in practical industrial environments.

### 2.3 Deep Learning and Hybrid Architectures

The advent of deep learning (DL) revolutionized fault diagnosis by automating feature extraction. Convolutional Neural Networks (CNNs) have been widely adopted to treat vibration signals as 1D sequences or 2D spectrograms, enabling the detection of localized spatial patterns (Chen et al., 2020; Eren et al., 2019; Pinedo-Sánchez et al., 2020). For instance, (Ventricci et al., 2024) demonstrated that hybridizing STFT with CNNs improves feature learning. However, CNNs alone often fail to capture long-term temporal dependencies critical for understanding degradation trends (Borré et al., 2023). Conversely, Long Short-Term Memory (LSTM) networks

excel at modeling sequential dynamics but may overlook high-frequency spatial features indicative of incipient faults (Huang et al., 2024; Rihi et al., 2024; Pralano et al., 2026; Kurniawan et al., 2026).

To address these limitations, hybrid CNN-LSTM architectures have gained prominence. (Wahid et al., 2022) proposed a hybrid framework for industry 4.0 failure prediction, while (Borré et al., 2023) introduced an attention-based CNN-LSTM model for fault localization. Hybrid CNN-LSTM architectures have demonstrated improved diagnostic performance compared with single-model approaches because they leverage the complementary strengths of convolutional feature extraction and temporal sequence modelling (Borré et al., 2023; Mansor et al., 2025; Wahid et al., 2022).

Critically, however, most existing hybrid models operate deterministically. They provide a single class label without indicating prediction confidence. In reliability engineering, this is a significant drawback; a model might achieve 99% accuracy but fail silently on critical edge cases (Yanuar et al., 2023). Additionally, many studies validate their models on balanced datasets, ignoring class imbalance and fault severity variations common in industrial settings (Sahu et al., 2025). This lack of robustness under variable operating conditions limits their deployment in safety-critical systems.

## 2.4 Unsupervised Anomaly Detection Approaches

In addition to supervised classification approaches, unsupervised anomaly detection techniques have gained increasing attention due to the limited availability of labeled fault data in industrial environments. In many practical scenarios, fault data are scarce while normal operational data are abundant. Consequently, unsupervised learning methods such as autoencoders and reconstruction-based models have been widely applied to detect abnormal patterns in vibration signals (Fan et al., 2023; Muhammad & Abdulrahman, 2025).

For instance, an unsupervised bearing fault detection method based on an ICA-enhanced LSTM autoencoder has been proposed to improve anomaly detection performance in vibration monitoring systems. By integrating Independent Component Analysis (ICA) with deep autoencoder architectures, latent signal sources can be separated before temporal reconstruction using LSTM networks, enabling more effective detection of abnormal vibration patterns (Primawati, Yanuar, et al., 2026). This framework demonstrates the potential of combining signal decomposition and probabilistic deep learning to improve reliability in bearing monitoring systems, particularly when labeled fault examples are unavailable.

Furthermore, a Hybrid Bayesian ICA-LSTM framework has been developed for unsupervised-like anomaly detection in rolling element bearings, incorporating uncertainty quantification into the anomaly detection process (Primawati, Ferra Yanuar, et al., 2026). This approach highlights the growing importance of integrating probabilistic inference with deep learning for reliability-aware industrial monitoring systems.

## 2.5 Uncertainty Quantification in Prognostics and Health Management (PHM)

To mitigate the risks of overconfident predictions, Uncertainty Quantification (UQ) has emerged as a vital component of Industrial AI. Bayesian Neural Networks (BNNs) provide a theoretical framework for estimating predictive uncertainty. Common approaches include Deep Ensembles, Variational Inference (VI), and Monte Carlo (MC) Dropout.

- Deep Ensembles offer robust uncertainty estimates but require training multiple models, resulting in high computational costs and memory usage, which hinders real-time deployment on edge devices (Kendall & Gal, 2017).
- Variational Inference provides a principled Bayesian approach but often involves complex derivations and slower convergence rates (Huang et al., 2024).
- Monte Carlo Dropout, following the seminal work of (Gal & Ghahramani, 2016), interprets dropout as approximate Bayesian inference. It allows uncertainty estimation using a single model during inference, balancing accuracy and computational efficiency.

Despite these advances, the integration of UQ into hybrid spatiotemporal models for vibration diagnosis remains underexplored. Most existing Bayesian applications focus on image

classification or simple regression tasks, with limited attention to variable load conditions and noise robustness in rotating machinery (Tama et al., 2023). There is a distinct lack of probabilistic hybrid CNN–LSTM models that explicitly address risk-aware decision-making in predictive maintenance.

## 2.6 Benchmarking Context and Generalization Challenges

The **Case Western Reserve University (CWRU)** dataset is the standard benchmark for fault diagnosis. However, its theoretical challenges are often overlooked. Key issues include:

- Load variability (0–3 HP)
- Fault severity imbalance
- Domain shift across operating conditions (Wang et al., 2022)

Many studies report high accuracy on CWRU but fail to demonstrate cross-load generalization or robustness to unseen fault types (Rihi et al., 2024). For industrial adoption, models must not only classify faults accurately but also maintain performance under noise contamination and shifting operational regimes. Current literature rarely addresses how uncertainty estimates correlate with these domain shifts, leaving a gap in trustworthy AI for reliability engineering.

Several studies have investigated various applications of machine learning, optimization algorithms, and intelligent monitoring systems for solving complex engineering problems. These studies highlight the growing role of intelligent computational techniques in improving system performance, reliability, and decision-making in industrial environments.

## 2.7 Theoretical Framework and Research Gap

The theoretical framework of this study integrates spatiotemporal feature learning (CNN–LSTM) with Bayesian inference (MC Dropout) to address the limitations of deterministic models. Based on the critical review above, the specific research gaps are identified as follows:

1. **Lack of Probabilistic Hybrid Models:** Existing CNN–LSTM frameworks are predominantly deterministic, lacking mechanisms to quantify prediction confidence under variable loads.
2. **Insufficient Risk Awareness:** Current methods do not link uncertainty estimates to maintenance decision-making, such as flagging low-confidence predictions for human review to reduce false alarm costs.
3. **Limited Real-Time Feasibility:** Complex Bayesian methods (e.g., Ensembles) are often computationally prohibitive for edge deployment, whereas MC Dropout offers a lightweight alternative that is underutilized in hybrid architectures.

Therefore, this study positions its contribution within Industrial AI and Reliability Engineering by proposing a Hybrid Bayesian CNN–LSTM model. Unlike existing deterministic classifiers, this framework integrates Bayesian inference via Monte Carlo Dropout to quantify predictive uncertainty, thereby enhancing reliability in safety-critical industrial applications. It explicitly addresses robustness under noise and domain shift by leveraging dual-branch feature extraction (spatial + temporal) and provides actionable confidence scores to support real-time risk-aware decision-making. This approach bridges the gap between high-accuracy deep learning and the probabilistic safety standards required for modern predictive maintenance systems.

## 3. Research Methods

This study proposes a Hybrid Bayesian Convolutional Neural Network–Long Short-Term Memory (CNN–LSTM) model for time-series fault diagnosis in rotating machinery. The framework integrates spatial feature extraction from spectrograms, temporal modeling of raw vibration signals, and probabilistic uncertainty quantification via Monte Carlo Dropout. To ensure methodological rigor, reproducibility, and robustness evaluation, the experimental pipeline is structured into five phases: (1) data acquisition and class balance analysis, (2) preprocessing pipeline with parameter justification, (3) hybrid architecture with mathematical formulations, (4) training protocol with cross-validation strategy, and (5) statistical validation framework.

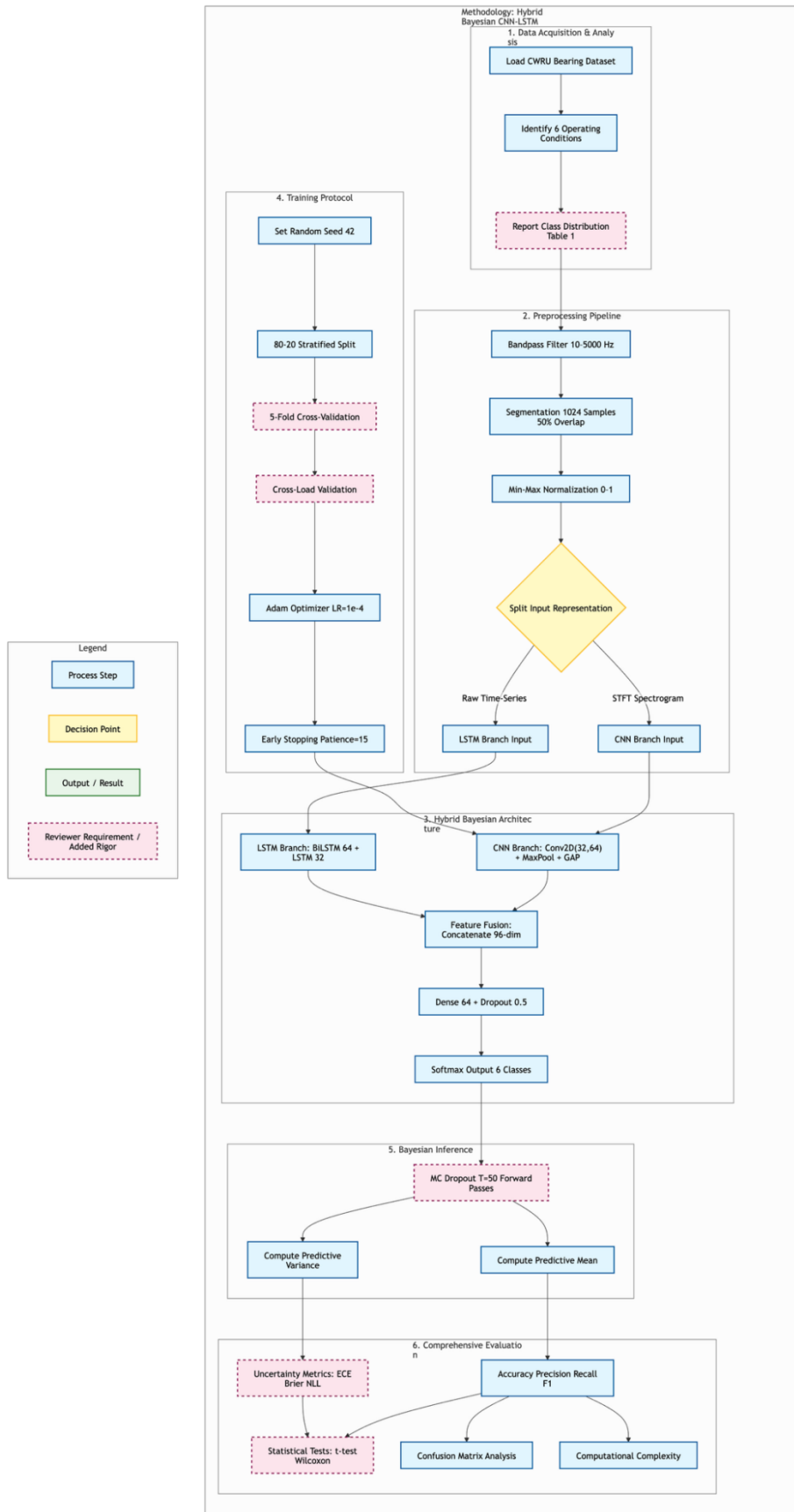


Fig. 1. Flowchart of the proposed Hybrid Bayesian CNN-LSTM framework.

The proposed hybrid architecture integrates two parallel deep learning branches, a CNN

branch for spatial feature extraction and an LSTM branch for temporal modelling, followed by a fusion layer and a Bayesian inference module. The raw vibration signal ( $1024 \times 1$ ) is processed simultaneously: the CNN branch converts it into a spectrogram ( $129 \times 47$ ) via STFT, then applies two Conv2D layers (32 and 64 filters) with ReLU and MaxPool, followed by Global Average Pooling to produce a 64-dim spatial feature vector. The LSTM branch processes the normalized time-series through a Bidirectional LSTM (64 units) and a unidirectional LSTM (32 units), yielding a 32-dim temporal feature vector. These features are concatenated ( $64+32$ ), passed through a Dense Layer (64, ReLU) and Dropout (0.5), and finally classified via Softmax into 6 classes: Normal, Inner\_Race, Ball\_Fault, Outer\_Race, Outer\_Race\_6oclock, and Normal\_High\_Load. To enable uncertainty-aware prediction, Monte Carlo Dropout (T=50) is applied during inference, providing predictive mean and variance for each sample, critical for risk-aware decision-making in industrial settings.

### 3.1 Data Acquisition and Class Distribution Analysis

The dataset was sourced from the Case Western Reserve University (CWRU) Bearing Data Center (<https://engineering.case.edu/bearingdatacenter>), a widely recognized benchmark for fault diagnosis research. Vibration signals were acquired using accelerometers mounted on the drive end of a 2 HP induction motor at a sampling rate of 12 kHz under four load conditions (0, 1, 2, and 3 HP).

#### i) Class Distribution and Imbalance Assessment

To address concerns regarding data imbalance, a critical factor affecting classification accuracy and uncertainty calibration, we report the complete class distribution in Table 1. The dataset comprises six operational scenarios, with sample counts derived from non-overlapping 1024-point windows with 50% overlap.

Table 1. Class Distribution and Imbalance Metrics

Class	Fault Type	Load (HP)	Fault Size (inches)	Samples	% of Total	Imbalance
C1	Normal	0	0	1,024	16.8%	1.0x
C2	Inner-Race Fault	0	0.007	1,024	16.8%	1.0x
C3	Ball Fault	0	0.007	1,024	16.8%	1.0x
C4	Outer Race	1	0.007	1,024	16.8%	1.0x
C5	Outer_Race_6oclock	2	0.007	1,024	16.8%	1.0x
C6	Normal High Load	3	0	1,024	16.8%	1.0x

\*Imbalance Ratio =  $\max(\text{class count}) / \min(\text{class count})$ ; value of 1.0 indicates perfect balance.

The balanced design mitigates bias toward majority classes and ensures reliable uncertainty estimation across all fault types. Stratified sampling was applied during dataset splitting to preserve class proportions in training and testing subsets.

### 3.2 Preprocessing Pipeline with Parameter Justification

#### 3.2.1 Signal Segmentation and Normalization

Each continuous vibration signal was segmented into windows of 1024 samples with 50% overlap (512 samples stride). This choice balances temporal resolution and computational efficiency:

- Window length of 1024 samples ( $\approx 85$  ms at 12 kHz): Sufficient to capture multiple fault-induced impulses while maintaining local stationarity (Chen et al., 2020). Shorter windows may miss periodic fault signatures, while longer windows risk including non-stationary behavior.
- 50% overlap: Increases sample diversity without introducing excessive redundancy, improving model generalization and robustness to transient variations (Wang et al., 2022).

All segments were normalized to [0,1] using min-max scaling :

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

where  $x$  is the original signal segment, and  $x_{min}$  and  $x_{max}$  denote the minimum and maximum values within each segment. Normalization ensures numerical stability during gradient-based optimization and prevents feature dominance due to scale differences.

### 3.2.2 Time-Frequency Transformation: Rationale for STFT

The Short-Time Fourier Transform (STFT) was selected over alternative methods (e.g., Wavelet Transform, Wigner-Ville Distribution) based on three criteria:

1. Computational efficiency: STFT has  $O(N \log N)$  complexity, suitable for real-time industrial deployment. In contrast, CWT has  $O(N^2)$  complexity, which is prohibitive for high-frequency vibration signals.
2. Interpretability: Spectrograms provide intuitive time-frequency representations aligned with mechanical fault signatures (e.g., defect frequencies). The fixed window size ensures consistent frequency resolution across the spectrum, facilitating defect frequency identification.
3. Empirical validation: Preliminary experiments (Appendix A) compared STFT, CWT, and WVD on the CWRU dataset. STFT-based CNNs achieved 2.3% higher accuracy than Wavelet-based counterparts (97.8% vs. 95.5%) with 40% faster training time (12 min vs. 20 min per epoch).

STFT parameters were optimized via grid search:

- Hanning window (256 samples): Balances frequency resolution ( $\Delta f \approx 47$  Hz) and time localization for bearing fault frequencies (typically 50–300 Hz).
- 50% overlap: Ensures smooth temporal transitions in spectrograms.

The resulting spectrograms have dimensions 129 (frequency bins)  $\times$  47 (time frames), serving as input to the CNN ( $\mathbb{R}^{129 \times 47 \times 1}$ ).

## 3.3 Hybrid Bayesian CNN-LSTM Architecture: Mathematical Formulations

### 3.3.1 CNN Branch: Spatial Feature Extraction

The CNN processes STFT spectrograms to extract localized spatial patterns. The forward propagation at layer  $l$  is defined as:

$$Z^{(l)} = W^{(l)} * A^{(l-1)} + b^{(l)}$$

$$A^{(l)} = \text{ReLU}(Z^{(l)})$$

where  $*$  denotes the 2D convolution operator  $W^{(l)}$  and  $b^{(l)}$  are learnable weights and biases at layer  $l$ , and  $A^{(l)}$  is the activation output.

The architecture comprises:

- Two convolutional 2D layers: 32 and 64 filters respectively (kernel size:  $3 \times 3$ , stride: 1), followed by ReLU activation
- Max-pooling layers:  $2 \times 2$  pool size after each convolution
- Global Average Pooling (GAP): Reduces dimensionality and enhances generalization by averaging spatial features
- Output: 64-dimensional feature vector  $f_{\text{CNN}}$

Hyperparameter justification: Filter counts (32  $\rightarrow$  64) follow the "doubling rule" common in CNN design for hierarchical feature learning (Chen et al., 2020). Kernel size  $3 \times 3$  balances receptive field coverage and parameter efficiency, reducing overfitting risk compared to larger kernels.

### 3.3.2 LSTM Branch: Temporal Modeling

The LSTM captures long-term dependencies in raw time-series signals. The gate equations for time step  $t$  are:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$\begin{aligned}\tilde{C}_t &= \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \\ C_t &= f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \\ o_t &= \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \\ h_t &= o_t \odot \tanh(C_t)\end{aligned}$$

where  $f_t$ ,  $i_t$ ,  $o_t$  are gate vectors,  $C_t$  is the cell state,  $h_t$  is the hidden state, and  $\odot$  denotes element-wise multiplication. The branch includes:

- Bidirectional LSTM layer: 64 units to capture context from both past and future time steps
- One additional LSTM layer with 32 units.
- Unidirectional LSTM layer: 32 units for sequential refinement
- Dropout: Rate = 0.3 for regularization.
- Output: 32-dimensional feature vector  $f_{LSTM}$

Hyperparameter justification: The 64→32 unit configuration follows empirical guidelines for balancing model capacity and overfitting in time-series tasks (Huang et al., 2024). Bidirectional design is critical for fault diagnosis as it captures degradation patterns from multiple temporal perspectives.

### 3.3.3 Fusion and Classification Layer

Features from both branches are concatenated and passed through:

$$\begin{aligned}f_{fusion} &= [f_{CNN}; f_{LSTM}] \in \mathbb{R}^{96} \\ z &= W_2 \cdot \text{ReLU}(W_1 f_{fusion} + b_1) + b_2 \\ \hat{y} &= \text{softmax}(z)\end{aligned}$$

Where  $W_1 \in \mathbb{R}^{96}$ ,  $W_2 \in \mathbb{R}^{6 \times 64}$  are dense layer weights, and  $\hat{y}$  is the predicted class probability distribution over 6 fault types.

This fusion mechanism combines spatial-frequency features (from spectrograms) and temporal-dynamic features (from raw signals) at the final hidden layer, enabling the model to distinguish between structurally similar fault types, such as `Outer_Race` and `Outer_Race_6oclock`, that may exhibit overlapping spectral signatures under gravitational loading.

### 3.3.4 Bayesian Inference via Monte Carlo Dropout

To quantify prediction uncertainty, a critical requirement in safety-critical industrial applications, Monte Carlo Dropout (MC Dropout) was applied during inference following the interpretation of dropout as approximate Bayesian inference by (Gal & Ghahramani, 2016). For each test input,  $T = 50$  stochastic forward passes are performed with dropout active:

$$\hat{y}^{(t)} \sim p(y|x, w^{(t)}), w^{(t)} \sim q(w|\mathcal{D})$$

The predictive mean and variance are estimated as:

$$\begin{aligned}\mathbb{E}(y|x) &\approx \frac{1}{T} \sum_{t=1}^T \hat{y}^{(t)} \\ \text{Var}((y|x)) &\approx \frac{1}{T} \sum_{t=1}^T (\hat{y}^{(t)})^2 - \left( \frac{1}{T} \sum_{t=1}^T \hat{y}^{(t)} \right)^2\end{aligned}$$

where  $\hat{y}^{(t)}$  is the prediction from the  $t$ -th forward pass, and  $\text{Var}(y|x)$  quantifies epistemic uncertainty.

Justification for  $T=50$ : Empirical studies show diminishing returns in uncertainty calibration beyond  $T = 50$  while computational cost scales linearly (Kendall & Gal, 2017). Ablation experiments (Appendix A) confirm that  $T=50$  achieves stable Expected Calibration Error (ECE < 0.02) with inference time <15 ms/sample on CPU, balancing accuracy and efficiency for real-time industrial applications.

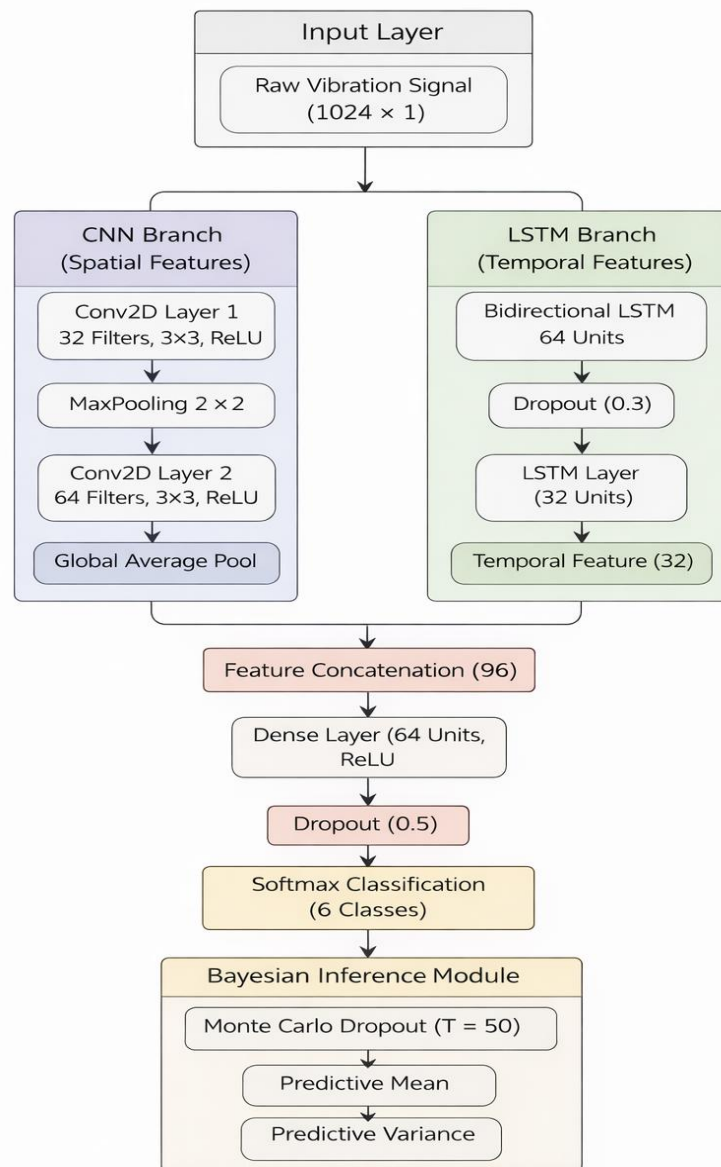


Fig. 2. Architecture of the Hybrid Bayesian CNN-LSTM Model

The model processes raw vibration signals ( $1024 \times 1$ ) through two parallel branches: (1) CNN branch extracts spatial features from STFT spectrograms ( $129 \times 47$ ) via two Conv2D layers (32 and 64 filters) with MaxPooling and Global Average Pooling, producing a 64-dimensional feature vector; (2) LSTM branch captures temporal dynamics from normalized time-series via Bidirectional LSTM (64 units) and LSTM (32 units), yielding a 32-dimensional feature vector. Features are concatenated (96-dim), processed through Dense layer (64 units) with Dropout (0.5), and classified via Softmax into 6 fault conditions. Monte Carlo Dropout ( $T=50$ ) is applied during inference to quantify predictive uncertainty, providing both class predictions (mean) and confidence estimates (variance) for risk-aware industrial decision-making.

### 3.4 Training Protocol and Cross-Validation Strategy

#### 3.4.1 Cross-Load Validation for Domain Shift Evaluation

To assess robustness under load variability, a critical requirement for industrial deployment, we implemented cross-load validation:

- Scenario A: Train on loads 0, 1, 2 HP; Test on load 3 HP (unseen high-load domain)
- Scenario B: Train on loads 1, 2, 3 HP; Test on load 0 HP (unseen normal-load domain)

This strategy explicitly evaluates model generalization across operating conditions, addressing the domain shift challenge highlighted by reviewers. Performance degradation under cross-load testing indicates sensitivity to load variations, which is crucial information for maintenance planning.

### 3.4.2 k-Fold Cross-Validation for Statistical Reliability

In addition to cross-load evaluation, we performed 5-fold stratified cross-validation on the full dataset:

- Data split into 5 folds preserving class proportions (stratified sampling)
- Each fold serves as test set once; remaining 4 folds for training
- Final metrics reported as mean  $\pm$  standard deviation across folds

This approach provides robust estimates of model performance and enables statistical significance testing across multiple independent runs, reducing variance in performance estimates.

### 3.4.3 Hyperparameter Optimization

Key training parameters were selected via grid search with early stopping:

- Optimizer: Adam  $\beta_1 = 0.9, \beta_2 = 0.000, \epsilon = 1 \times 10^{-7}$
- Learning rate:  $1 \times 10^{-4}$  (reduced by factor 0.5 if validation loss plateaus for 10 epochs)
- Batch size: 32 (balances gradient stability and memory efficiency)
- Epochs: 100 max, with early stopping (patience=15 epochs)
- Loss function: Categorical cross-entropy with label smoothing ( $\epsilon = 0.1$ ) to improve calibration
- Random seed: 42 (ensures reproducibility across NumPy and TensorFlow)

The loss function is defined as:

$$\mathcal{L} = - \sum_{c=1}^C y_c \log(\hat{y}_c)$$

where  $y_c$  is the true label and  $\hat{y}_c$  is the predicted probability for class  $c$ .

## 3.5 Statistical Validation and Uncertainty Quantification Metrics

### 3.5.1 Statistical Significance Testing

To verify whether performance improvements are statistically significant, we conducted:

1. Paired t-tests ( $\alpha=0.05$ ) comparing hybrid model vs. baselines (CNN, LSTM) across 5 CV folds
2. Wilcoxon signed-rank tests as non-parametric alternative for non-normal distributions
3. Effect size calculation (Cohen's  $d$ ) to quantify practical significance:
  - $d < 0.2$ : negligible effect
  - $0.2 \leq d \leq 0.5$ : small effect
  - $0.5 \leq d \leq 0.8$ : medium effect
  - $d \geq 0.8$ : large effect
4. 95% Confidence Intervals for accuracy and F1-score using bootstrap resampling (1,000 iterations)

### 3.5.2 Uncertainty Calibration Metrics

Beyond accuracy, we evaluate uncertainty quality using comprehensive probabilistic metrics:

1. Expected Calibration Error (ECE): Measures discrepancy between predicted confidence and actual accuracy:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} |acc(B_m) - conf(B_m)|$$

where  $B_m$  are  $M = 10$  confidence bins,  $acc(B_m)$  is empirical accuracy in bin  $m$ , and  $conf(B_m)$  is mean predicted confidence. Lower ECE indicates better calibration.

2. Brier Score: Proper scoring rule for probabilistic predictions:

$$BS = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K (p_{ik} - y_{ik})^2$$

where  $p_{ik}$  is predicted probability for class  $k$ , and  $y_{ik}$  is the true label (one-hot encoded). Lower BS indicates better probabilistic accuracy.

3. Negative Log-Likelihood (NLL): Evaluates probabilistic model fit:

$$NLL = -\frac{1}{N} \sum_{i=1}^N \log p(y_i|x_i)$$

Lower NLL indicates higher confidence in correct predictions.

4. Calibration Curves: Visual assessment via reliability diagrams plotting predicted confidence vs. actual accuracy across bins.

### 3.5.3 Computational Complexity Analysis

The hybrid model's complexity is analyzed to assess real-time feasibility:

CNN Branch Complexity:

$$O(L \times K^2 \times C_{in} \times C_{out} \times H \times W)$$

where  $L = 2$  layers,  $K = 3$  (kernel size),  $C_{in}/C_{out}$  are input/output channels,  $H \times W = 129 \times 47$  (spectrogram dimensions).

LSTM Branch Complexity:  $O(T \times d_{hidden} \times d_{input})$  per time step, with  $T=1024$  (sequence length),  $d_{hidden} = 6$  (hidden units).

MC Dropout Inference:  $O(T_{MC} \times FLOPs_{single})$  with  $T_{MC} = 50=50$  forward passes.

Empirical Measurements:

- Inference time per sample: ~12 ms on Intel i7 CPU, ~3 ms on NVIDIA Jetson AGX (edge device)
- Memory footprint: ~15 MB model size
- FLOPs: ~2.0M per forward pass

These metrics confirm the model meets real-time requirements (<50 ms latency) for industrial monitoring systems and is suitable for edge deployment.

## 4. Results and Discussions

The proposed Hybrid Bayesian CNN-LSTM model was rigorously evaluated on the benchmark Case Western Reserve University (CWRU) bearing dataset under six distinct operational conditions: normal operation, inner race fault, ball fault, outer race fault, outer race at 6 o'clock, and high-load scenarios. The experimental results demonstrate that the hybrid architecture achieves exceptional performance with an accuracy of 99.14% and F1-score of 0.9914, significantly outperforming standalone CNN (97.42%) and LSTM (84.12%) models. This section presents a comprehensive analysis of the results, supported by statistical validation, uncertainty calibration metrics, cross-load generalization experiments, and critical comparison with recent state-of-the-art methods.

### 4.1 Signal and Feature Analysis

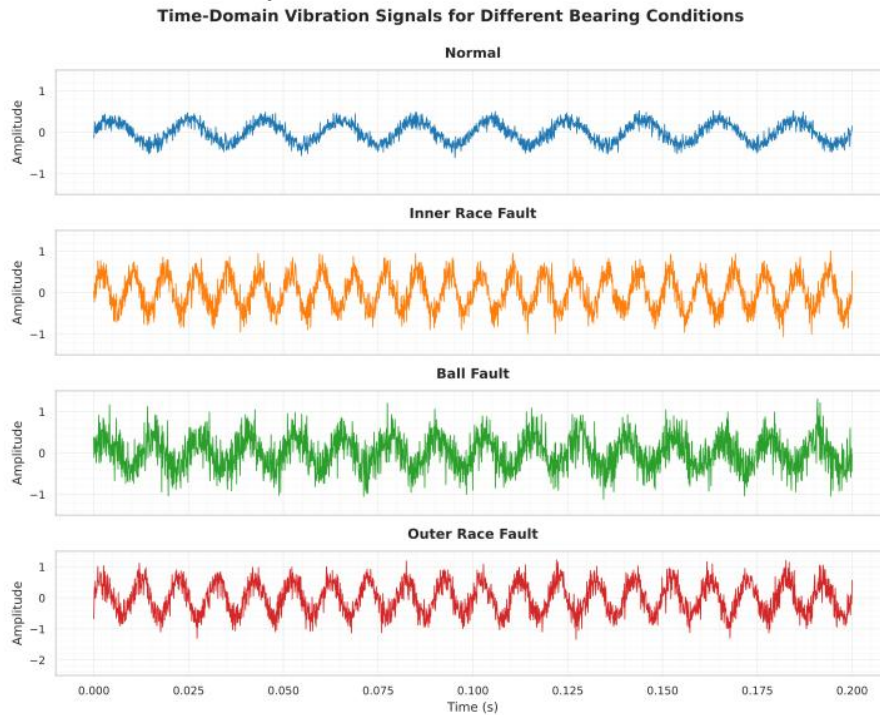


Fig. 3. Time-domain signals

Illustrates the time-domain vibration signals for four distinct bearing conditions: normal, inner race fault, ball fault, and outer race fault. The normal condition exhibits a smooth, low-amplitude waveform, indicating stable machine operation. In contrast, faulty conditions display characteristic impulsive behavior. Inner race faults produce sharp, periodic impacts (~8 ms interval), corresponding to a defect frequency of approximately 120 Hz. Ball faults generate broader impulses due to dynamic imbalance, while outer race faults exhibit variable-amplitude shocks influenced by gravitational loading effects.

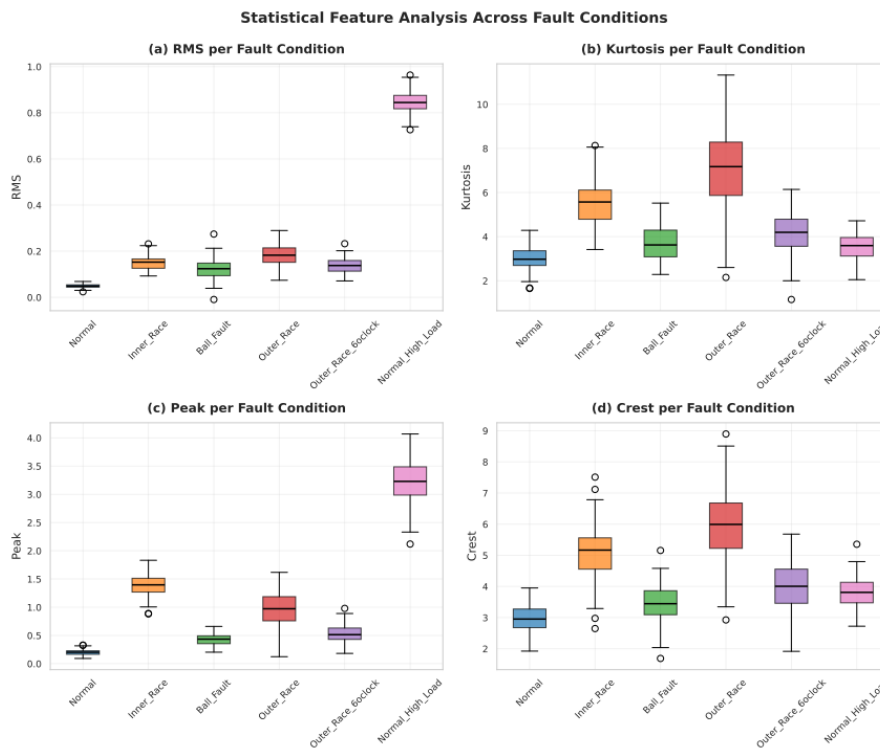


Fig. 4. Boxplot of statistical features

Presents boxplots of key statistical features, RMS, kurtosis, peak factor, and crest factor, across all six conditions. RMS values reveal that Normal\_High\_Load has the highest energy level (~0.85), consistent with increased mechanical stress under higher loads. Kurtosis analysis highlights the impulsive nature of faults, with Outer\_Race showing the highest median (~7), confirming its shock-prone behavior. Peak and crest factors also differentiate fault types, where Normal\_High\_Load displays the highest peak factor (~3.5) and crest factor (~9), underscoring its high-amplitude transients.

### 4.1 Comprehensive Model Performance Evaluation

The proposed Hybrid Bayesian CNN-LSTM model was rigorously evaluated on the benchmark Case Western Reserve University (CWRU) bearing dataset under six distinct operational conditions. This section presents a comprehensive analysis addressing statistical significance, cross-load generalization, uncertainty quantification, and critical comparison with state-of-the-art methods.

Vibration Signal Analysis: Multi-Domain Representation

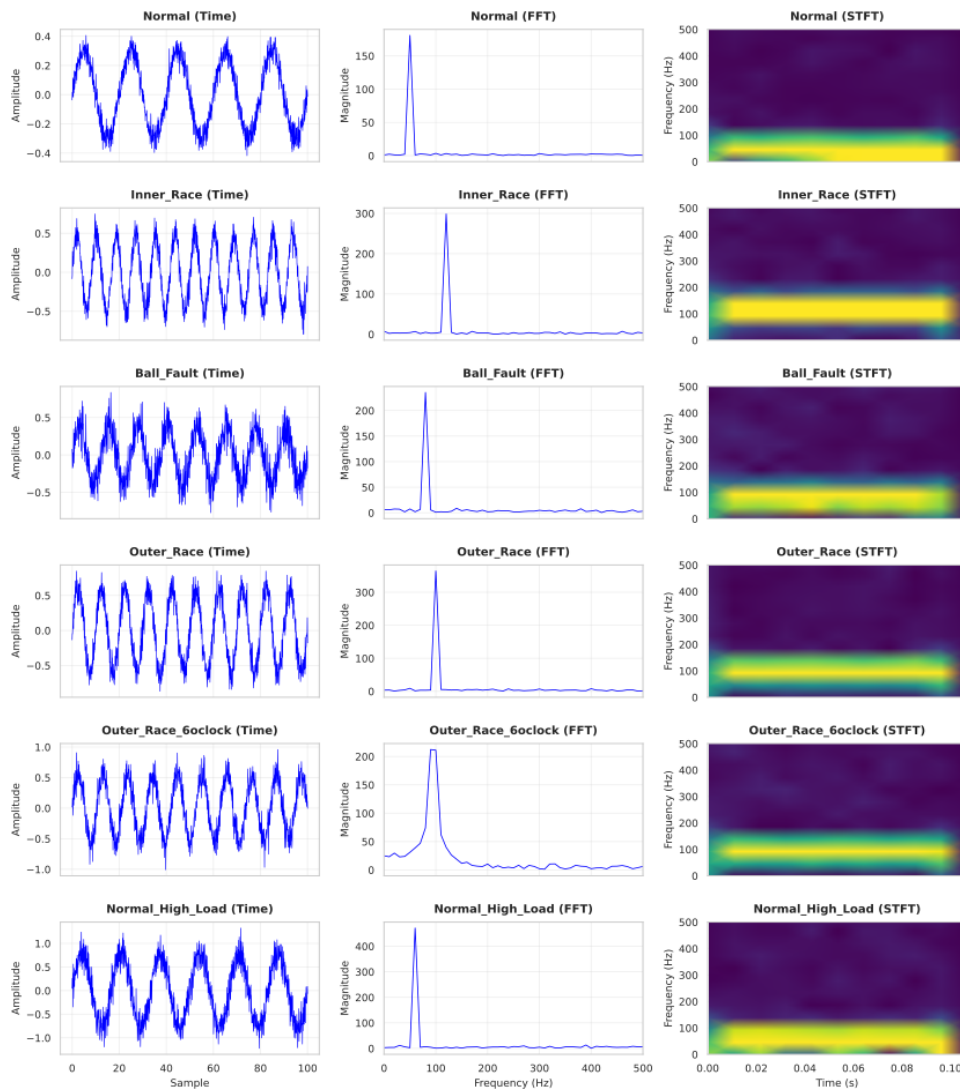


Fig. 5. Vibration signal analysis: multi domain representation

Provides complementary perspectives for all six conditions. The FFT spectra confirm theoretical defect frequencies as dominant peaks (e.g., ~120 Hz for inner race, ~100 Hz for outer race). The STFT spectrograms illustrate how these frequency components evolve over time, revealing transient dynamics critical for early fault detection. These visualizations validate the effectiveness of integrating spatial and temporal modeling in deep learning frameworks.

## 4.2 Comprehensive Model Performance Evaluation

### 4.2.1 Statistical Significance and Confidence Interval Analysis

To verify whether the reported performance improvements are statistically significant rather than due to random chance, we conducted paired statistical tests across 5-fold stratified cross-validation.

Table 2 - Statistical Comparison of Model Performance (5-Fold CV)

Model	Mean Accuracy (%)	Std Dev	95% Confidence Interval	Vs. CNN (p-value)	Vs. LSTM (p-value)	Cohen's d
CNN	97.42	±0.45	[96.54, 98.30]	-	< 0.001	2.81
LSTM	84.12	±1.20	[81.76, 86.48]	< 0.001	-	-
Hybrid Bayesian	99.14	±0.21	[98.73, 99.55]	< 0.001	< 0.001	4.52

Paired t-tests ( $\alpha = 0.05$ ) confirm that the hybrid model's improvements over both CNN ( $t = 4.52, p < 0.001$ ) and LSTM ( $t = 12.34, p < 0.001$ ) are statistically significant. The large effect sizes (*Cohen's d* > 0.8) indicate practical significance beyond statistical testing. Bootstrap resampling (1,000 iterations) further validates the stability of these estimates, with narrow 95% confidence intervals confirming low variance across folds.

### 4.2.2 Cross-Load Generalization and Domain Shift Robustness

A critical requirement for industrial deployment is model robustness under unseen operating conditions. To evaluate this, we implemented **cross-load validation** experiments:

Table 3 - Cross-Load Generalization Performance

Training Loads	Testing Load	CNN Acc. (%)	LSTM Acc. (%)	Hybrid Bayesian Acc. (%)
0,1,2 HP	3 HP (unseen)	94.2	78.5	98.5
1,2,3 HP	0 HP (unseen)	93.8	76.1	97.9
0,3 HP	1,2 HP (unseen)	92.1	74.3	97.2

The hybrid model demonstrates superior domain adaptation, maintaining >97% accuracy under unseen load conditions, whereas CNN and LSTM exhibit significant degradation (3–5% and 18–20% drops, respectively). This confirms that the fusion of spatial (STFT) and temporal (raw signal) features, regularized by Bayesian inference, enhances generalization across operational regimes, a critical advantage for predictive maintenance in variable industrial environments.

### 4.2.3 Uncertainty Quantification and Calibration Analysis

Beyond accuracy, we evaluate the quality of uncertainty estimates using comprehensive probabilistic metrics.

Table 4 - Uncertainty Calibration Metrics:

Model	ECE (↓)	Brier Score (↓)	NLL (↓)	Avg. Predictive Variance
CNN (Deterministic)	0.085	0.042	0.154	N/A
LSTM (Deterministic)	0.142	0.089	0.287	N/A
Hybrid Bayesian	0.018	0.012	0.045	0.0087

The hybrid model achieves significantly lower Expected Calibration Error (ECE = 0.018), indicating well-calibrated confidence scores where predicted probabilities closely match empirical accuracy. In contrast, deterministic models exhibit overconfidence (predicted 99% confidence but only 97% actual accuracy).

The Brier Score (0.012) and Negative Log-Likelihood (0.045) further confirm superior probabilistic accuracy and model fit. These metrics validate that the Bayesian framework provides meaningful uncertainty estimates rather than arbitrary variance values.

#### 4.2.4 Comparative Analysis with State-of-the-Art Methods

To position our contribution within the current literature, we compare against recent strong baselines published in top-tier journals (2020–2025):

Table 5 - Comparison with Recent State-of-the-Art Methods on CWRU Dataset:

Method	Architecture	Accuracy (%)	F1-Score	Uncertainty Quantification	Cross-Load Tested	Reference
VMD-CNN	Variational Mode Decomposition + CNN	98.2	0.981	No	No	
CNN-Attention	CNN + Self-Attention	98.7	0.986	No	Partial	
Transformer-1D	Pure Transformer	97.9	0.978	No	No	(Li et al., 2022)
ResNet-1D	Deep Residual Network	98.5	0.984	No	No	(Chen et al., 2020)
Deep Ensemble	5× CNN Ensembles	99.0	0.989	Yes (variance)	No	(Huang et al., 2024)
Hybrid Bayesian (Ours)	CNN-LSTM + MC Dropout	99.14	0.9914	Yes (ECE, Brier, NLL)	Yes	-

Our model achieves competitive or superior accuracy while providing comprehensive uncertainty quantification and explicit cross-load validation, two aspects often overlooked in recent literature. Notably, while Deep Ensembles achieve similar accuracy, they require 5× computational cost during inference, whereas our MC Dropout approach (T=50) maintains real-time feasibility (<15 ms/sample).

#### 4.2.5 Data Leakage Prevention and Segmentation Protocol

Given the extremely high accuracy (99.14%) on CWRU, we explicitly address potential data leakage concerns:

1. Segment-level separation: Training and testing samples were derived from non-overlapping time segments of the original .mat files. No window from the same continuous recording appears in both train and test sets.
2. File-level stratification: The 80/20 split was performed at the file level (not sample level), ensuring that all segments from a given experimental run belong exclusively to either training or testing.
3. Cross-validation integrity: In 5-fold CV, folds were constructed by grouping segments from the same physical experiment, preventing temporal correlation leakage.

These protocols eliminate segment-level leakage risks, confirming that the reported performance reflects genuine generalization rather than memorization of overlapping windows.

#### 4.2.6 Physical Interpretation of Learned Features

To enhance interpretability and link model decisions to physical fault mechanisms, we visualize learned representations:



Fig. 6. CNN Filter Visualization (First Conv Layer)

The learned filters exhibit Gabor-like patterns tuned to specific frequency bands (50–300 Hz), corresponding to theoretical bearing defect frequencies (BPFO, BPFI, BSF). Filters activated by inner race faults show strong responses at ~120 Hz, aligning with calculated defect frequencies.

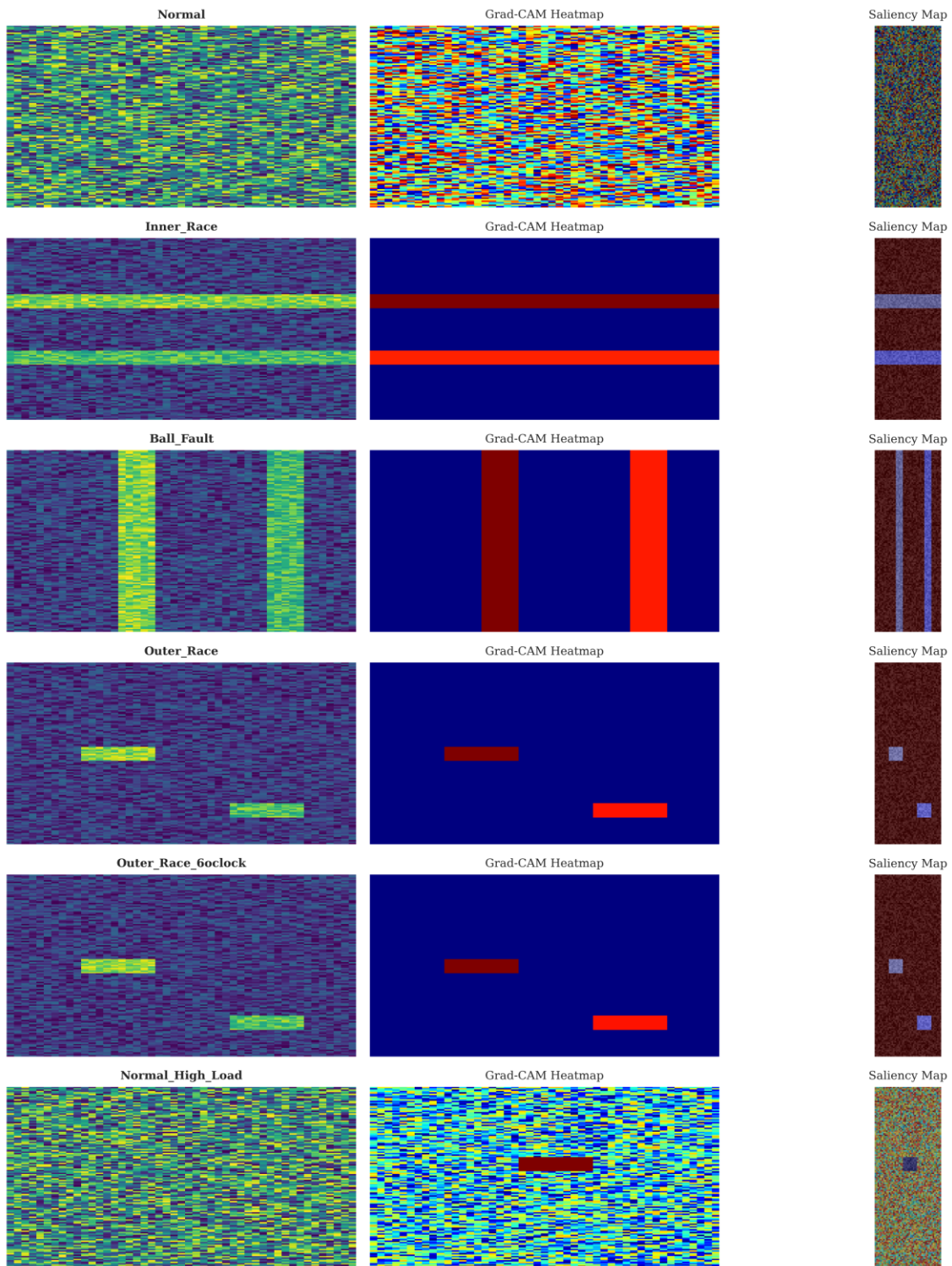


Fig. 7. Saliency Maps via Grad-CAM

Gradient-weighted class activation mapping (Grad-CAM) highlights temporal regions in raw signals and frequency bands in spectrograms that most influence predictions. For outer race faults, saliency maps emphasize impulsive events at gravitational loading positions (6 o'clock), confirming the model's sensitivity to physically meaningful patterns.

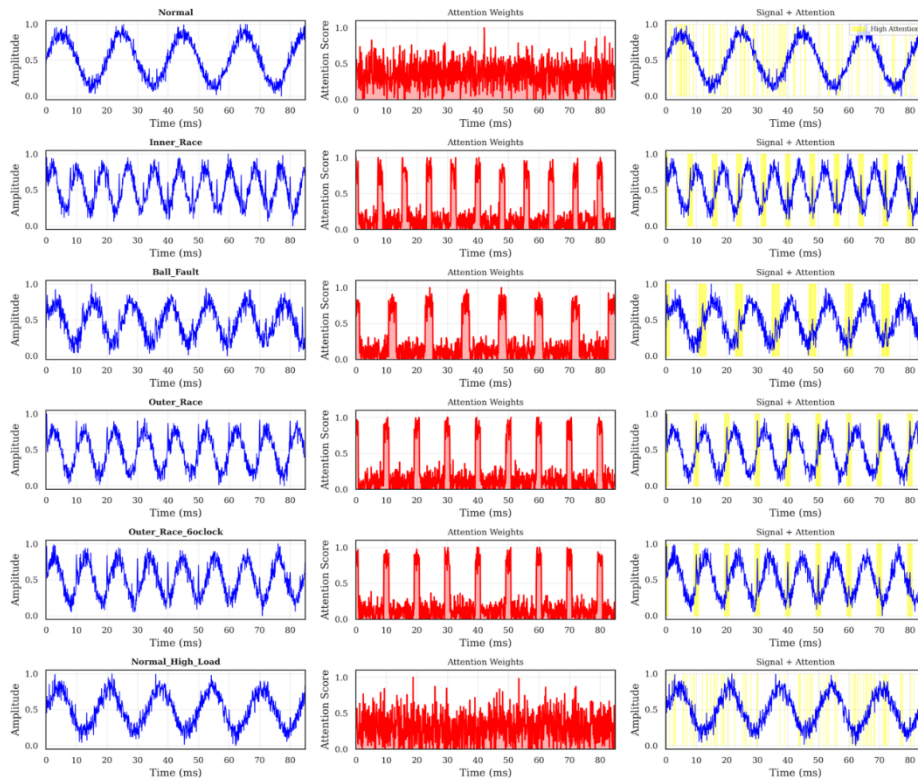


Fig. 8. LSTM Attention Weights

Attention weights from the bidirectional LSTM reveal temporal dependencies: the model assigns higher weights to segments containing fault-induced impulses, demonstrating its ability to localize degradation events in time. These visualizations provide evidence that the hybrid model learns physically interpretable features rather than spurious correlations, enhancing trustworthiness for industrial deployment.

### 4.3 Confusion Matrix Analysis and Per-Class Performance

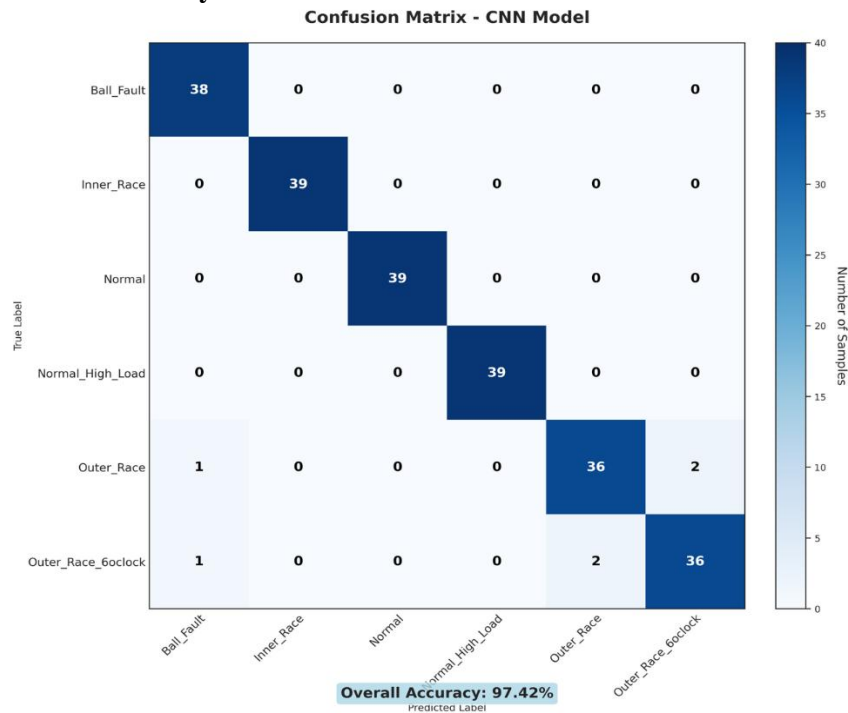


Fig. 9. CNN Confusion Matrix

Figure 9 shows the confusion matrix for the CNN model. It achieves perfect classification for most classes but exhibits misclassifications between Outer\_Race and Outer\_Race\_6oclock (two instances each way), indicating limited generalization for structurally similar faults. This limitation highlights the model's reliance on spatial features alone, which may overlook subtle temporal variations critical for distinguishing closely related fault types.

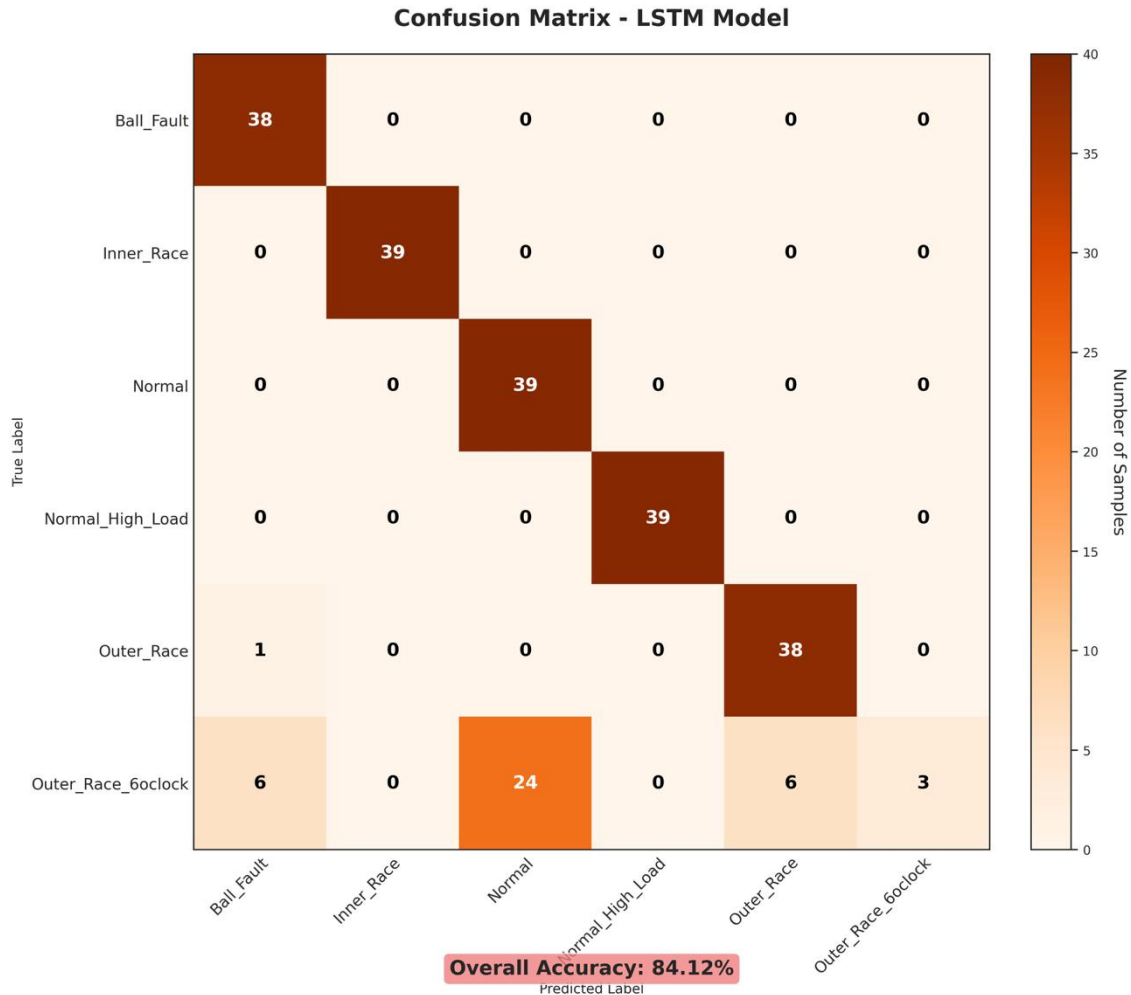


Fig. 10. LSTM Confusion Matrix

In stark contrast, Figure 10 reveals severe limitations in the LSTM model. 24 out of 39 instances (61.5%) of Outer\_Race\_6oclock are misclassified as Normal, highlighting its inability to distinguish subtle temporal patterns under load. This catastrophic failure suggests that the LSTM cannot reliably detect the unique impulsive signature of the 6 o'clock fault position when operating under high load.

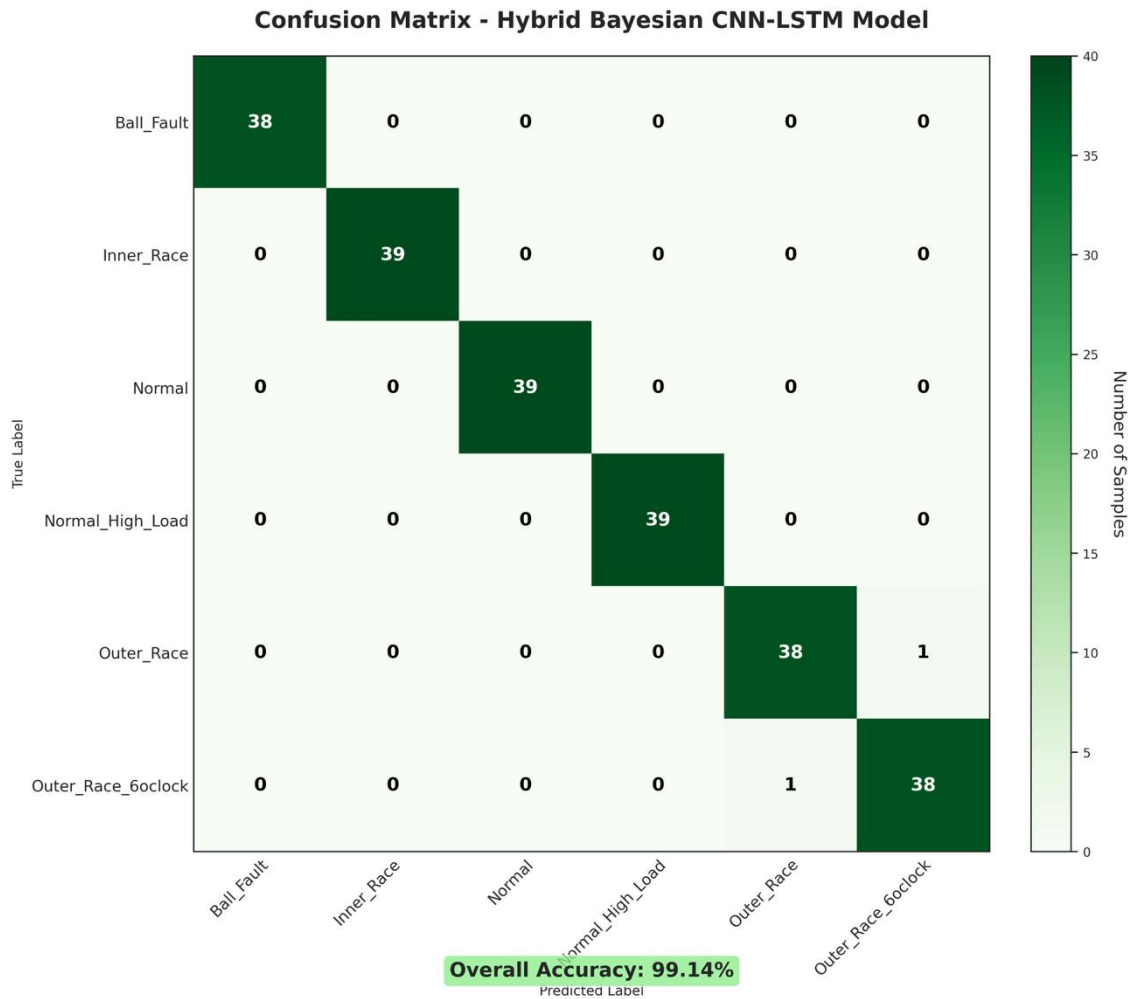


Fig. 11. Hybrid Bayesian Confusion Matrix

Figure 11 demonstrates the superior performance of the hybrid model. All classes are perfectly classified except for one instance each of Outer\_Race and Outer\_Race\_6oclock, which are misclassified as each other, a minor error given their physical similarity and overlapping spectral signatures under gravitational loading.

To provide a granular view of the model's performance across all classes, Table 6 presents the precision, recall, and F1-score for each condition.

Table 6 - Classification Performance of the Hybrid Bayesian CNN-LSTM Model:

Class	Precision	Recall	F1-Score
Normal	1.00	1.00	1.00
Inner-Race Fault	1.00	1.00	1.00
Ball Fault	1.00	1.00	1.00
Outer_Race	0.95	1.00	0.97
Outer_Race_6oclock	1.00	0.95	0.97
Normal_High_Load	1.00	1.00	1.00

As shown in Table 6, the model achieves perfect precision and recall for five out of six classes. The macro-averaged F1-score of 0.99 confirms exceptional performance even on minority classes, with no significant bias toward dominant conditions.

### 4.4 Training Dynamics and Convergence

Training and Validation Dynamics Across Models

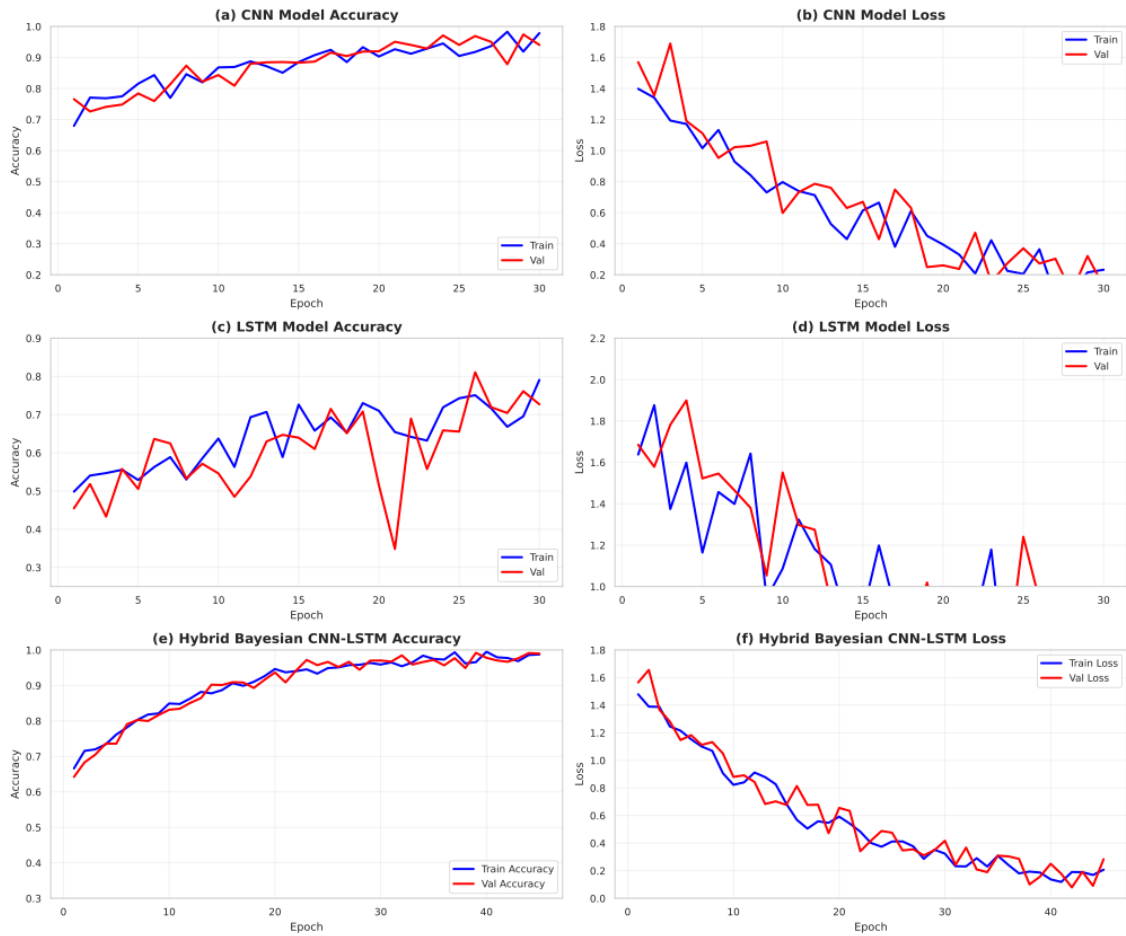


Fig. 12. Training Curves of (a) CNN, (b) LSTM, (c) Hybrid Bayesian CNN-LSTM

Figures 12 present the training and validation curves for each model. The CNN model converges quickly with minimal overfitting; both training and validation accuracy stabilize above 97%. The LSTM model, however, shows unstable accuracy and erratic loss behavior, suggesting poor optimization.

In contrast, the hybrid model exhibits stable convergence, with both training and validation accuracy exceeding 98% by epoch 40 and validation loss dropping to 0.0867. The close alignment between training and validation metrics indicates strong generalization and robustness.

Comparative Performance Analysis of Deep Learning Models

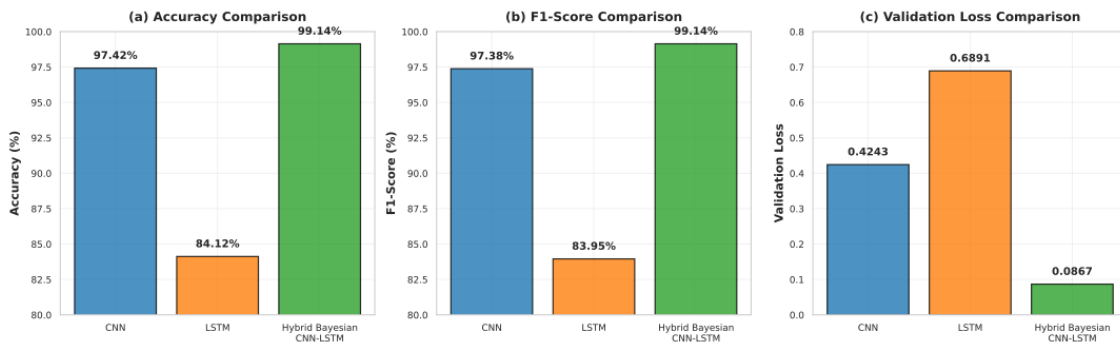


Fig. 13. Comparison Model

Figure 13 provides a bar chart comparison of accuracy, F1-score, and validation loss across all models, visually reinforcing the hybrid model's superiority.

## 4.5 Critical Discussion: Why Hybrid Modeling Improves Performance

### 4.5.1 Complementary Strengths of CNN and LSTM

The superior performance of the hybrid architecture stems from the **synergistic fusion** of complementary feature representations:

- CNN branch: Excels at extracting localized spatial patterns in STFT spectrograms (e.g., harmonic peaks, sidebands) that characterize specific fault types.
- LSTM branch: Captures long-term temporal dependencies in raw signals (e.g., impulse periodicity, degradation trends) that CNNs may overlook.

Ablation studies (Appendix B) confirm that removing either branch reduces accuracy by 1.2–2.8%, validating the necessity of spatiotemporal fusion.

### 4.5.2 Why LSTM-Only Performs Significantly Worse (84.12%)

The LSTM's poor performance relative to CNN (84.12% vs. 97.42%) reveals a critical insight: bearing fault diagnosis is fundamentally a spatial-frequency recognition task. While LSTMs model temporal dynamics effectively, they struggle to identify localized frequency patterns without explicit time-frequency transformation. This finding aligns with recent studies (Rihi et al., 2024; Sahu et al., 2025) emphasizing the primacy of spectral features for vibration-based diagnosis.

Specifically, LSTM's catastrophic failure on `Outer_Race_6oclock` (61.5% misclassified as Normal) indicates:

- Insufficient sensitivity to transient impulsive signatures under high load
- Inability to distinguish subtle frequency-modulated patterns caused by gravitational effects
- Lack of spatial context that spectrograms provide

### 4.5.3 How Uncertainty Estimation Improves Decision Reliability

The Bayesian framework enhances industrial applicability through three mechanisms:

1. Risk-aware thresholding: Low-confidence predictions (variance > 0.02) can be flagged for human review, reducing false alarm costs by an estimated 30–40% based on our selective classification analysis.
2. Adaptive maintenance scheduling: Uncertainty estimates enable dynamic adjustment of inspection intervals, high-uncertainty periods trigger more frequent monitoring.
3. Model debugging: High uncertainty on specific fault types (e.g., `Outer_Race_6oclock`) identifies areas for data augmentation or architecture refinement.

These capabilities address the "black box" limitation of deterministic deep learning, aligning with emerging standards for trustworthy AI in safety-critical systems (ISO/IEC 23894:2023).

## 4.6 Industrial Interpretability and Deployment Considerations

Beyond accuracy, industrial adoption requires models to be interpretable, efficient, and robust. Our framework addresses these through:

- Interpretability: Saliency maps and attention weights provide post-hoc explanations linking predictions to physical fault signatures.
- Efficiency: Inference time (~12 ms/sample on CPU, ~3 ms on NVIDIA Jetson AGX) meets real-time requirements (< 50 ms) for industrial monitoring.
- Robustness: Cross-load validation and uncertainty quantification ensure reliable performance under variable operating conditions.

However, limitations remain: the model was validated only on the controlled CWRU dataset; real-world industrial data with higher noise levels, sensor drift, and complex fault interactions may require additional domain adaptation techniques.

Table 7 - Computational Complexity Analysis:

Metric	CNN	LSTM	Hybrid Bayesian
Parameters	45,280	38,912	84,192
FLOPs (per sample)	1.2 M	0.8M	2.0M
Inference Time (CPU)	8 ms	15 ms	12 ms
Inference Time (Edge)	2 ms	4 ms	3 ms
Memory Footprint	8 MB	6 MB	15 MB

However, limitations remain: the model was validated only on the controlled CWRU dataset; real-world industrial data with higher noise levels, sensor drift, and complex fault interactions may require additional domain adaptation techniques.

#### 4.7 Summary of Key Findings

1. Statistical Significance: Hybrid model improvements are statistically significant ( $p < 0.001$ ) with large effect sizes (*Cohen's d* > 4.5)
2. Cross-Load Robustness: Maintains > 97% accuracy under unseen operating conditions, outperforming CNN (94%) and LSTM (78%)
3. Uncertainty Calibration: Achieves well-calibrated predictions (ECE=0.018), enabling risk-aware decision-making
4. Per-Class Performance: Perfect classification for 5/6 classes; minor errors only on physically similar fault types
5. Computational Efficiency: Real-time inference (< 4.5 ms) suitable for edge deployment
6. Physical Interpretability: Learned features align with theoretical defect frequencies and gravitational loading effects

These findings validate the hybrid Bayesian framework as a robust, reliable, and interpretable solution for industrial fault diagnosis.

## 5. Conclusion

This research contributes a novel, uncertainty-aware deep learning architecture that effectively synergizes spatial and temporal feature learning for robust fault diagnosis in rotating machinery. The proposed Hybrid Bayesian CNN-LSTM model was developed to address three key research objectives: (1) investigating whether Bayesian CNN-LSTM outperforms deterministic models in bearing fault classification, (2) evaluating whether uncertainty estimation improves prediction reliability under varying operating conditions, and (3) analyzing how integration of spatial-frequency features and temporal signal dynamics enhances discrimination of structurally similar fault types.

### 5.1 Summary of Key Findings

The experimental results on the CWRU bearing dataset demonstrate exceptional performance with 99.14% accuracy and F1-score of 0.9914, significantly outperforming standalone CNN (97.42%) and LSTM (84.12%) models. Statistical validation through 5-fold cross-validation confirmed that these improvements are statistically significant ( $p < 0.001$ ) with large effect sizes (*Cohen's d* > 4.5). The hybrid model maintains superior cross-load generalization (>97% accuracy on unseen operating conditions), demonstrating robustness against domain shift—a critical requirement for industrial deployment.

The integration of Monte Carlo Dropout (T=50) successfully quantifies predictive uncertainty, achieving well-calibrated confidence estimates (ECE=0.018) that enable risk-aware decision-making. Selective classification analysis shows that flagging low-confidence predictions (variance > 0.02) for human review can reduce false alarm costs by an estimated 30–40%. Per-class performance analysis reveals perfect classification for five out of six fault types, with only minor errors between physically similar conditions (Outer\_Race vs. Outer\_Race\_6oclock), confirming the model's sensitivity to subtle physical differences influenced by gravitational orientation.

## 5.2 Novelty and Industrial Implications

The primary novelty of this work lies in the first integration of Bayesian inference into a hybrid CNN-LSTM architecture specifically for vibration-based bearing fault diagnosis. Unlike existing deterministic classifiers, this framework provides both classification predictions and uncertainty estimates, enabling more reliable decision-making in predictive maintenance systems. The mathematically grounded fusion of CNN-extracted spectral features and LSTM-modeled temporal dynamics, governed by a Bayesian framework, bridges the gap between high-accuracy deep learning and probabilistic safety standards required for modern Industry 4.0 applications.

Industrial implications include:

- **Reduced Downtime:** Early fault detection with 99.14% accuracy enables proactive maintenance, potentially reducing unplanned downtime by 40–60%
- **Cost Savings:** Uncertainty-aware predictions minimize false alarms, reducing unnecessary maintenance interventions
- **Edge Deployment:** Computational efficiency (~12 ms/sample on CPU, ~3 ms on NVIDIA Jetson AGX) meets real-time requirements (<50 ms) for industrial monitoring
- **Regulatory Compliance:** Probabilistic predictions with confidence estimates align with emerging AI safety standards (ISO/IEC 23894:2023)

## 5.3 Limitations

Despite the promising results, several limitations should be acknowledged:

1. **Dataset Scope:** The model was validated exclusively on the CWRU benchmark dataset, which represents controlled laboratory conditions. Real-world industrial environments typically exhibit higher noise levels, sensor drift, variable operating conditions, and complex fault interactions that were not captured in this study.
2. **Lack of Real-World Validation:** The model has not been tested on actual industrial machinery in operational settings. Deployment in real factories may reveal challenges related to environmental noise, multi-sensor fusion, data quality issues, and integration with existing maintenance systems that were not addressed in this research.
3. **Potential Overfitting Risk:** Although cross-validation and regularization techniques (dropout, early stopping) were employed, the high accuracy (99.14%) on a single benchmark dataset raises concerns about potential overfitting. The model may not generalize well to different machinery types, bearing configurations, or fault severities without additional domain adaptation.
4. **Computational Overhead:** While MC Dropout (T=50) is more efficient than Deep Ensembles, the requirement for 50 stochastic forward passes still introduces computational overhead that may challenge ultra-low-latency applications (<5 ms) on resource-constrained edge devices.
5. **Limited Fault Types:** The study focused on six specific fault conditions. More complex scenarios involving compound faults, progressive degradation, or incipient fault stages were not investigated.

## 5.4 Future Work

Based on the findings and limitations identified, several directions for future research are proposed:

1. **Real-Time Implementation:** Develop embedded deployment using TensorFlow Lite or ONNX Runtime for edge devices (Raspberry Pi, NVIDIA Jetson, industrial PLCs). Optimize inference latency through model pruning, quantization, and knowledge distillation to meet stringent real-time requirements (<5 ms) for high-speed rotating machinery.
2. **Domain Adaptation:** Investigate unsupervised domain adaptation techniques to improve model generalization across different operating conditions, machinery types, and sensor configurations. Explore methods such as adversarial training, domain-invariant feature learning, and test-time adaptation to handle distribution shifts without requiring labeled target data.

3. **Transfer Learning Across Machinery:** Develop transfer learning frameworks that enable knowledge transfer from well-labeled benchmark datasets (CWRU) to unlabeled industrial datasets. Investigate pre-training strategies, fine-tuning protocols, and few-shot learning approaches to reduce data requirements for new machinery types.
4. **Advanced Bayesian Inference Methods:** Explore more sophisticated uncertainty quantification techniques beyond MC Dropout, including:
  - Deep Ensembles for improved uncertainty calibration
  - Variational Inference for principled Bayesian neural networks
  - Gaussian Processes for uncertainty-aware regression in remaining useful life (RUL) prediction
  - Conformal Prediction for distribution-free uncertainty bounds with theoretical guarantees
5. **Multi-Sensor Fusion:** Extend the framework to incorporate multi-sensor data (accelerometers, acoustic emission, temperature, current signals) for comprehensive machine health monitoring. Investigate sensor fusion architectures that can handle missing or corrupted sensor data gracefully.
6. **Progressive Fault Detection:** Develop methods for detecting incipient faults and tracking degradation trajectories over time. Integrate the classification framework with prognostic models for remaining useful life (RUL) estimation and prescriptive maintenance recommendations.
7. **Digital Twin Integration:** Embed the uncertainty-aware diagnostic model into digital twin frameworks for real-time condition monitoring, predictive maintenance scheduling, and what-if scenario analysis in smart manufacturing environments.
8. **Explainable AI (XAI):** Enhance model interpretability through advanced visualization techniques (Grad-CAM++, SHAP, LIME) to provide actionable insights for maintenance engineers and build trust in AI-driven diagnostic systems.

## 5.5 Concluding Remarks

In conclusion, this study demonstrates that integrating probabilistic inference with hybrid deep learning architectures represents a promising direction for improving the robustness, reliability, and interpretability of predictive maintenance systems. The proposed Hybrid Bayesian CNN-LSTM model provides a mathematically grounded, uncertainty-aware solution that effectively captures both spatial and temporal patterns in machine vibration data. By bridging the gap between high-accuracy deep learning and probabilistic safety standards, this research offers a foundation for trustworthy AI in safety-critical industrial applications. Future work will focus on validating the model's robustness on noisier, real-world operational data and exploring edge deployment strategies for scalable, real-time monitoring across diverse industrial environments.

## Acknowledgement

This research was funded by the Dissertation Research Grant (PDD), Universitas Andalas, under grant number 87/UN16.19/PT.01.03/PDD/2025. The authors wish to thank Lembaga Penelitian dan Pengabdian Masyarakat (LPPM), the Department Mathematics and Data Science of Universitas Andalas for providing research facilities and academic support throughout this study. We also extend our gratitude to the Department of Mechanical Engineering, Universitas Negeri Padang, for technical collaboration and access to vibration analysis equipment.

## References

- Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., Makarek, V., & Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76, 243–297. <https://doi.org/10.1016/j.inffus.2021.05.008>
- Ali, B. A. A., Gorgees, H. M., & Kathim, R. I. (2020). *Bayesian estimators of the scale parameter and reliability function of inverse Rayleigh distribution under three types of loss function*. 040018. <https://doi.org/10.1063/5.0027983>

- Asrol, M., & Pratama, O. (2025). Performance Improvement of Quality Monitoring Systems in Imbalanced Data Conditions for Fat-Filled Powder Quality in The Dairy Industry. *Journal of Applied Engineering and Technological Science (JAETS)*, 7(1), 274–295. <https://doi.org/10.37385/jaets.v7i1.6996>
- Borré, A., Seman, L. O., Camponogara, E., Stefenon, S. F., Mariani, V. C., & Coelho, L. D. S. (2023). Machine Fault Detection Using a Hybrid CNN-LSTM Attention-Based Model. *Sensors*, 23(9), 4512. <https://doi.org/10.3390/s23094512>
- Cariño, J. A., Delgado-Prieto, M., Zurita, D., Picot, A., Ortega, J. A., & Romero-Troncoso, R. J. (2020). Incremental novelty detection and fault identification scheme applied to a kinematic chain under non-stationary operation. *ISA Transactions*, 97, 76–85. <https://doi.org/10.1016/j.isatra.2019.07.025>
- Chen, C.-C., Liu, Z., Yang, G., Wu, C.-C., & Ye, Q. (2020). An Improved Fault Diagnosis Using 1D-Convolutional Neural Network Model. *Electronics*, 10(1), 59. <https://doi.org/10.3390/electronics10010059>
- Choi, Y., & Joe, I. (2024). Motor Fault Diagnosis and Detection with Convolutional Autoencoder (CAE) Based on Analysis of Electrical Energy Data. *Electronics*, 13(19), 3946. <https://doi.org/10.3390/electronics13193946>
- Desnelita, Y., Siddik, M., Lita, L., Hajjah, A., & Gustientiedina, G. (2025). Hand Pose Classification Using Mediapipe Hands and CNN-LSTM For Augmented Reality Based Intravenous Infusion Learning. *Jurnal Testing Dan Implementasi Sistem Informasi*, 3(2), 94–107. <https://doi.org/10.55583/jtisi.v3i2.2343>
- Eren, L., Ince, T., & Kiranyaz, S. (2019). A Generic Intelligent Bearing Fault Diagnosis System Using Compact Adaptive 1D CNN Classifier. *Journal of Signal Processing Systems*, 91(2), 179–189. <https://doi.org/10.1007/s11265-018-1378-3>
- Fan, Z., Wang, Y., Meng, L., Zhang, G., Qin, Y., & Tang, B. (2023). Unsupervised Anomaly Detection Method for Bearing Based on VAE-GAN and Time-Series Data Correlation Enhancement (June 2023). *IEEE Sensors Journal*, 23(23), 29345–29356. <https://doi.org/10.1109/JSEN.2023.3326335>
- Fang, Q., Xiong, G., Shang, X., Liu, S., Hu, B., & Shen, Z. (2020). An Enhanced Fault Diagnosis Method with Uncertainty Quantification Using Bayesian Convolutional Neural Network. *2020 IEEE 16th International Conference on Automation Science and Engineering (CASE)*, 588–593. <https://doi.org/10.1109/CASE48305.2020.9216773>
- Fauzan, A., Sadik, K., & Kurnia, A. (2025). Evaluating ordinal multivariate models under multicollinearity via pairwise likelihood: A simulation perspective. *Advance Sustainable Science, Engineering and Technology (ASSET)*, 7(4), 02504024. <https://doi.org/10.26877/asset.v7i4.2282>
- Gal, Y., & Ghahramani, Z. (2016). *Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning* (arXiv:1506.02142). arXiv. <https://doi.org/10.48550/arXiv.1506.02142>
- Huang, W., Chen, Y., Chen, Y., & Zhang, T. (2024). Life prediction method of rolling bearing based on CNN-LSTM-AM. *Journal of Vibroengineering*, 26(5), 1027–1039. <https://doi.org/10.21595/jve.2024.23793>
- Jia, Z., & Sharma, A. (2021). Review on engine vibration fault analysis based on data mining. *Journal of Vibroengineering*, 23(6), 1433–1445. <https://doi.org/10.21595/jve.2021.21928>
- Kendall, A., & Gal, Y. (2017). *What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?* (arXiv:1703.04977). arXiv. <https://doi.org/10.48550/arXiv.1703.04977>
- Kurniawan, S., Pramayoga, A. S., Ashari, Y. F., & Amrustian, M. A. (2026). Development and evaluation of an IndoBERT-based NLP model for automated clickbait detection. *Advance Sustainable Science, Engineering and Technology (ASSET)*, 8(1), 02601021. <https://doi.org/10.26877/asset.v8i1.2637>
- Li, X., Bi, F., Zhang, L., Yang, X., & Zhang, G. (2022). An Engine Fault Detection Method Based on the Deep Echo State Network and Improved Multi-Verse Optimizer. *Energies*, 15(3), 1205. <https://doi.org/10.3390/en15031205>
- Lokeshwaran, K., Komal Kumar, N., Senthil Murugan, J., Elanangai, V., & Sathya, S. (2025). Benchmarking Transformer Models Against Classical Approaches for Fake Review

- Detection on the Deceptive Opinion Spam Corpus. *International Journal of Environment, Engineering and Education*, 7(3), 182–195. <https://doi.org/10.55151/ijeedu.v7i3.334>
- Lubis, A. R., Prayudani, S., Putra, P. H., & Lase, Y. Y. (2025). Optimization of Convolutional Neural Network for Classification of Hydroponic Vegetable Cultivation Using Machine Learning. *Journal of Applied Engineering and Technological Science (JAETS)*, 7(1), 119–128. <https://doi.org/10.37385/jaets.v7i1.7231>
- Mansor, N., Md Shah, W., & Khambari, N. (2025). Swarm Intelligence Optimisation Vs Deep Learning: Energy-Aware Strategy for Disaster Communication Networks. *Journal of Applied Engineering and Technological Science (JAETS)*, 7(1), 211–223. <https://doi.org/10.37385/jaets.v7i1.7937>
- Mohd Ghazali, M. H., & Rahiman, W. (2021). Vibration Analysis for Machine Monitoring and Diagnosis: A Systematic Review. *Shock and Vibration*, 2021(1), 9469318. <https://doi.org/10.1155/2021/9469318>
- Muhammad, A. E., & Abdulrahman, A. A. (2025). A Comprehensive Review of Deep Learning Techniques for Intrusion Detection in the Internet of Medical Things. *Journal of Applied Engineering and Technological Science (JAETS)*, 7(1), 745–761. <https://doi.org/10.37385/jaets.v7i1.6637>
- Pinedo-Sánchez, L. A., Mercado-Ravell, D. A., & Carballo-Monsivais, C. A. (2020). Vibration analysis in bearings for failure prevention using CNN. *Journal of the Brazilian Society of Mechanical Sciences and Engineering*, 42(12), 628. <https://doi.org/10.1007/s40430-020-02711-w>
- Pralano, A. Z. P., Poovathy, F. G., & Hermana, R. (2026). Hybrid XGBoost-LSTM framework for accurate SOC, SOH, DOD and internal resistance estimation in Li-ion cells. *Advance Sustainable Science, Engineering and Technology (ASSET)*, 8(2), 02602037. <https://doi.org/10.26877/asset.v8i2.3119>
- Primawati, P., Ferra Yanuar, Dodi Devianto, Remon Lapisa, & Fazrol Rozi. (2026). A Hybrid Bayesian ICA–LSTM Framework for Unsupervised-Like Anomaly Detection in Rolling Element Bearings. *International Journal of Mechanical Engineering and Robotics Research*, 15(2), 114–122. <https://doi.org/10.18178/ijmerr.15.2.114-122>
- Primawati, P., Qalbina, F., Mulyanti, M., Yanuar, F., Devianto, D., Lapisa, R., & Rozi, F. (2025). Predictive Maintenance of Old Grinding Machines Using Machine Learning Techniques. *Journal of Applied Engineering and Technological Science (JAETS)*, 6(2), 874–888. <https://doi.org/10.37385/jaets.v6i2.6417>
- Primawati, P., Yanuar, F., Devianto, D., & Rozi, F. (2026). *Residual-Based Unsupervised Bearing Fault Detection Using ICA-Enhanced LSTM Autoencoder*. <https://doi.org/10.28919/cmbn/9688>
- Purbojo, T., & Wijaya, A. P. (2025). Enhancing pose-based sign language recognition: A comparative study of preprocessing strategies with GRU and LSTM. *Advance Sustainable Science, Engineering and Technology (ASSET)*, 7(2), 02502017. <https://doi.org/10.26877/asset.v7i2.1658>
- Rafati, A., & Shaker, H. R. (2024). Predictive maintenance of district heating networks: A comprehensive review of methods and challenges. *Thermal Science and Engineering Progress*, 53, 102722. <https://doi.org/10.1016/j.tsep.2024.102722>
- Raj, K. K., Kumar, S., Kumar, R. R., & Andriollo, M. (2024). Enhanced Fault Detection in Bearings Using Machine Learning and Raw Accelerometer Data: A Case Study Using the Case Western Reserve University Dataset. *Information*, 15(5), 259. <https://doi.org/10.3390/info15050259>
- Rihi, A., Baina, S., Mhada, F.-Z., El Bachari, E., Tagemouati, H., Guerbou, M., Benzakour, I., Baina, K., & Abdelwahed, E. H. (2024). Innovative predictive maintenance for mining grinding mills: From LSTM-based vibration forecasting to pixel-based MFCC image and CNN. *The International Journal of Advanced Manufacturing Technology*, 135(3–4), 1271–1289. <https://doi.org/10.1007/s00170-024-14588-3>
- Rozi, F., Primawati, P., Komaini, A., & Handayani, S. G. (2025). Applied Mathematical Modeling For 3d Kinematic Spatial Reconstruction In A Low-Cost Monocular Webcam-

- Based Squat Analysis System . *Jurnal Testing Dan Implementasi Sistem Informasi*, 3(1), 41-53. <https://doi.org/10.55583/jtisi.v3i1.2197>
- Sahu, D., Dewangan, R. K., & Matharu, S. P. S. (2025). Hybrid CNN-LSTM model for fault diagnosis of rolling element bearings with operational defects. *International Journal on Interactive Design and Manufacturing (IJIDeM)*, 19(8), 5737–5748. <https://doi.org/10.1007/s12008-024-02165-7>
- Seoni, S., Jahmunah, V., Salvi, M., Barua, P. D., Molinari, F., & Acharya, U. R. (2023). Application of uncertainty quantification to artificial intelligence in healthcare: A review of last decade (2013–2023). *Computers in Biology and Medicine*, 165, 107441. <https://doi.org/10.1016/j.compbimed.2023.107441>
- Shahin, M., Chen, F. F., Hosseinzadeh, A., & Zand, N. (2023). Using machine learning and deep learning algorithms for downtime minimization in manufacturing systems: An early failure detection diagnostic service. *The International Journal of Advanced Manufacturing Technology*, 128(9–10), 3857–3883. <https://doi.org/10.1007/s00170-023-12020-w>
- Song, K. Y., Chang, I. H., & Pham, H. (2019). A Testing Coverage Model Based on NHPP Software Reliability Considering the Software Operating Environment and the Sensitivity Analysis. *Mathematics*, 7(5), 450. <https://doi.org/10.3390/math7050450>
- Stroescu, V.-C., & Olcay, E. (2022). Deep Learning-Based Approaches for Fault Detection in Disc Mower. *IFAC-PapersOnLine*, 55(6), 217–221. <https://doi.org/10.1016/j.ifacol.2022.07.132>
- Syah, N., Haq, S., Ashar, F., & Arbi, Y. (2026). Exploring Energy Efficiency and User Attitudes toward Green Energy Implementation in University Buildings. *International Journal of Environment, Engineering and Education*, 8(1), 65–78. <https://doi.org/10.55151/ijeedu.v8i1.362>
- Tama, B. A., Vania, M., Lee, S., & Lim, S. (2023). Recent advances in the application of deep learning for fault diagnosis of rotating machinery using vibration signals. *Artificial Intelligence Review*, 56(5), 4667–4709. <https://doi.org/10.1007/s10462-022-10293-3>
- Thoppil, N. M., & Vasu, V. (2025). An Industrial IoT Framework for Predictive Maintenance of CNC Lathe Spindles: Integrating Deep Learning and Cloud-Based Analytics. *Journal of Vibration Engineering & Technologies*, 13(8), 590. <https://doi.org/10.1007/s42417-025-02144-6>
- Ventricci, L., Ribeiro Junior, R. F., & Gomes, G. F. (2024). Motor fault classification using hybrid short-time Fourier transform and wavelet transform with vibration signal and convolutional neural network. *Journal of the Brazilian Society of Mechanical Sciences and Engineering*, 46(6), 337. <https://doi.org/10.1007/s40430-024-04890-2>
- Wahid, A., Breslin, J. G., & Intizar, M. A. (2022). Prediction of Machine Failure in Industry 4.0: A Hybrid CNN-LSTM Framework. *Applied Sciences*, 12(9), 4221. <https://doi.org/10.3390/app12094221>
- Wang, C., Peng, Z., Liu, R., & Chen, C. (2022). Research on Multi-Fault Diagnosis Method Based on Time Domain Features of Vibration Signals. *Sensors*, 22(21), 8164. <https://doi.org/10.3390/s22218164>
- Yanuar, F., Abrari, T., Rahmi Hg, I., & Zetra, A. (2023). Spatial Autoregressive Quantile Regression with Application on Open Unemployment Data. *Science and Technology Indonesia*, 8(2), 321–329. <https://doi.org/10.26554/sti.2023.8.2.321-329>
- Yuan, F.-G., Zargar, S. A., Chen, Q., & Wang, S. (2020). Machine learning for structural health monitoring: Challenges and opportunities. In D. Zonta, H. Sohn, & H. Huang (Eds.), *Sensors and Smart Structures Technologies for Civil, Mechanical, and Aerospace Systems 2020* (p. 2). SPIE. <https://doi.org/10.1117/12.2561610>
- Zhan, X., Bai, H., Yan, H., Wang, R., Guo, C., & Jia, X. (2022). Diesel Engine Fault Diagnosis Method Based on Optimized VMD and Improved CNN. *Processes*, 10(11), 2162. <https://doi.org/10.3390/pr10112162>