

Industrial Machine Fault Detection Using Hybrid CNN-LSTM Model with Real-Time SCADA Integration for Mechanical Fault Analysis

M. Priyadharshan, M.E(Ph.D), Assistant Professor

Arul Kumar, Sam Kevinadel, Sanjai

Department of Artificial Intelligence and Data Science Nehru Institute of Engineering and Technology

Thirumalayampalayam, Tamil

Mentor: Priyadharshan Sir

Nadu 641105, India

Abstract—This paper presents a comprehensive intelligent industrial machine fault detection system based on a hybrid Convolutional Neural Network (CNN) and Long Short-Term Memory (LSTM) architecture integrated with a real-time SCADA-style monitoring platform. The proposed system analyzes multivariate sensor data—vibration, temperature, pressure, current, and ambient air quality—to identify five categories of abnormal machine behavior with high precision. A Flask-based web dashboard with WebSocket communication provides real-time visualization of machine health status, live waveform plots, and a fault event log with confidence scores. A NodeMCU ESP8266 hardware controller activates a piezoelectric buzzer upon fault detection and displays status locally on an LCD module. The integration of deep learning inference with real-time hardware response enables proactive predictive maintenance, reducing unplanned downtime and improving industrial safety standards. The system is exhaustively validated on a three-phase induction motor testbed instrumented with six sensor modalities under five fault conditions: bearing failure, overheating, pressure anomaly, vibration excess, and electrical fault. The hybrid CNN-LSTM model achieves 97.4% overall classification accuracy, precision and recall exceeding 96% across all fault classes, a 1.8% false positive rate, and an end-to-end system latency of 87 milliseconds, significantly outperforming conventional threshold-based SCADA methods, standalone CNN, standalone LSTM, SVM with FFT features, Random Forest, and other contemporary baselines. This paper provides detailed treatment of the system architecture, dataset construction methodology, model design and training procedure, hardware integration, experimental results, and future research directions.

Index Terms—Industrial Fault Detection, CNN-LSTM, Predictive Maintenance, IoT Monitoring, Deep Learning, Flask Dashboard, NodeMCU ESP8266, Real-Time Monitoring, SCADA Integration, Multivariate Sensor Fusion, Bidirectional LSTM, Bearing Fault Detection, Industry 4.0, Edge Inference, TensorFlow Lite.

I. INTRODUCTION

Industrial machinery forms the cornerstone of modern manufacturing, chemical processing, power generation, and transportation infrastructure. The reliable and continuous operation of these machines directly determines production efficiency, product quality, operational safety, and overall profitability. Unexpected equipment failures not only cause production stoppages and downstream supply chain disruptions but also impose significant financial penalties and, in safety-critical environments, can endanger human lives. Conservative industry estimates place the average cost of unplanned manufacturing downtime at \$260,000 per hour [1], with large automotive or semiconductor fabrication plants reporting losses exceeding \$1 million per hour during critical production windows. These economic realities have made early, accurate, and automated fault detection one of the most intensely studied problems at the intersection of industrial engineering and applied machine learning.

The global industrial maintenance market has historically been dominated by two paradigms: corrective maintenance and time-based preventive maintenance. Corrective maintenance—repairing equipment only after failure—maximizes machine utilization but incurs the full economic cost of unexpected downtime, secondary damage, and emergency repair premiums. Time-based preventive maintenance replaces or services components on a fixed schedule regardless of their actual condition, leading to unnecessary expenditure on components with remaining useful life and, paradoxically, increased failure risk due to disturbances introduced during scheduled servicing. Both paradigms fail to leverage the rich continuous stream

of operational data generated by modern instrumented industrial machinery.

The emergence of condition-based maintenance (CBM) and, more recently, predictive maintenance (PdM) powered by machine learning represents a fundamental paradigm shift. Rather than responding to failures or adhering to fixed schedules, predictive maintenance continuously monitors equipment health indicators, learns the characteristic signature of normal operation, and provides advance warning of impending failures—enabling operators to plan interventions at the most cost-effective moment. The effectiveness of predictive maintenance is directly coupled to the quality of the underlying anomaly detection system.

Traditional automated fault detection has relied on SCADA (Supervisory Control and Data Acquisition) systems equipped with fixed threshold alarms. These systems trigger alerts when a monitored parameter exceeds a statically configured limit, such as a vibration amplitude threshold or a temperature setpoint. While providing a baseline level of protection, threshold-based systems suffer from fundamental limitations: they cannot adapt to changes in normal operating conditions caused by load variations, seasonal environmental changes, or gradual machine aging; they generate excessive false alarms under non-steady-state transients; and they cannot detect compound or incipient faults that manifest only through subtle multivariate interactions invisible to single-channel thresholds.

The rapid advancement of deep learning has created an opportunity to replace rigid threshold logic with adaptive, data-driven models that learn the complex multivariate manifold of normal machine operation and recognize deviations from it. Convolutional Neural Networks (CNNs) excel at extracting spatial features from raw sensor waveforms without manual feature engineering. Long Short-Term Memory (LSTM) networks and their bidirectional variants model temporal dependencies across sequential windows, capturing gradual degradation signatures that evolve over minutes or hours. The combination of these two architectures in a hybrid framework addresses both instantaneous anomalies and slow-developing fault patterns simultaneously.

A. Scope and Objectives

This project develops a complete, production-deployable predictive maintenance system for industrial rotating machinery. The specific objectives are: (1) design and train a hybrid CNN-LSTM deep learning model for multivariate industrial fault classification with at least 95% accuracy; (2) develop a real-time web monitoring dashboard using Flask and WebSocket technology providing sub-100 ms update latency; (3) implement a hardware alert layer using NodeMCU ESP8266 providing audible and visual fault notification within 50 ms of fault detection; (4) validate the complete integrated system on a three-phase induction motor testbed under five representative fault conditions; and (5) compare system performance against a comprehensive set of baseline methods including classical machine learning, standalone deep learning architectures, and conventional threshold-based SCADA.

B. Motivation and Problem Statement

Existing commercial SCADA installations at industrial facilities typically employ threshold-based alarms configured at commissioning time. These thresholds are determined from manufacturer specifications or initial operational experience and are rarely updated as machine condition evolves. As a result, modern industrial monitoring is characterized by a high false alarm burden—studies report false alarm rates of 10–40% in conventional vibration monitoring systems [11]—which leads to alarm fatigue and reduced operator responsiveness. Simultaneously, incipient faults that develop gradually below threshold levels are systematically missed until they produce catastrophic failure.

The core problem motivating this research is the absence of an integrated, low-cost system that combines the predictive accuracy of deep learning with the immediate physical response capability required by real industrial deployments. Most published deep learning fault detection research operates exclusively in the software domain, producing offline classification results without demonstrating practical integration with hardware alert systems or real-time monitoring interfaces. This paper addresses that gap by delivering a fully integrated hardware-software predictive maintenance solution validated under real operating conditions.

C. Contributions

The primary technical contributions of this work are: (i) a hybrid bidirectional CNN-LSTM architecture with five-class fault classification optimized for multivariate industrial sensor streams; (ii) a real-time Flask web application with WebSocket-based live

dashboard and persistent SQLite fault event logging; (iii) a complete NodeMCU ESP8266 hardware integration layer with severity-adaptive buzzer patterns and EEPROM fault logging; (iv) a six-channel multivariate sensor dataset collected from a three-phase induction motor testbed under controlled five-class fault conditions with 12,400 labeled samples; (v) comprehensive performance evaluation against eight baselines demonstrating consistent superiority across accuracy, precision, recall, F1-score, and false positive rate; and (vi) end-to-end latency profiling across all pipeline stages confirming production suitability.

D. Paper Organization

Section II provides an extended review of related work covering signal processing methods, deep learning architectures, and commercial SCADA systems. Section III details the proposed three-layer system architecture. Section IV presents the complete methodology including dataset construction, model architecture, and real-time inference pipeline. Section V describes each system module in depth. Section VI enumerates all hardware components and the complete software stack. Section VII presents experimental results including per-class metrics, latency profiling, ablation studies, and long-duration reliability tests. Section VIII provides detailed discussion of findings, limitations, and failure modes. Section IX outlines future research directions. Section X concludes the paper.

II. RELATED WORK AND EXISTING SYSTEMS

A. Signal Processing and Classical Machine Learning Approaches

The earliest automated fault detection methods in industrial settings applied statistical signal processing techniques to vibration signals captured by accelerometers mounted on rotating machinery. Fast Fourier Transform (FFT) analysis enabled frequency-domain identification of bearing defect frequencies, gear mesh frequencies, and rotor imbalance signatures. Envelope analysis of high-frequency resonance bands improved detectability of impulsive bearing faults. Wavelet decomposition provided time-frequency representations better suited to non-stationary fault signals than the purely frequency-domain FFT.

Features extracted by these signal processing techniques were subsequently classified using Support Vector Machines (SVM), achieving strong results on controlled laboratory datasets [2]. k -Nearest Neighbor (k NN), Naive Bayes, and Decision Tree classifiers were also evaluated extensively. While these methods demonstrated feasibility under controlled conditions, they required domain experts to select and engineer features for each new machine type, limiting their generalizability. Performance degraded significantly when operating conditions deviated from training distributions due to load changes, temperature variations, or machine aging.

B. Convolutional Neural Network Approaches

The emergence of deep learning provided a path to automatic feature extraction, eliminating the dependency on expert feature engineering. LeCun et al. [6] established the foundational principles of hierarchical convolutional feature learning. Zhang et al. [12] demonstrated that 1D CNNs applied directly to raw vibration waveforms outperformed hand-crafted FFT-SVM baselines on bearing fault detection tasks. Chen et al. [9] introduced transferable CNNs for rotary machinery fault diagnosis, demonstrating cross-domain generalization from laboratory testbeds to production environments.

Multi-scale CNN architectures applied parallel convolutional branches with different filter sizes to capture features across multiple frequency scales simultaneously, further improving accuracy on complex multi-fault scenarios.

C. Recurrent and Hybrid Architectures

Hochreiter and Schmidhuber [7] introduced Long Short-Term Memory networks, providing a solution to the vanishing gradient problem that had prevented effective training of recurrent networks on long sequences. LSTM networks modeled temporal dependencies within sequential sensor readings, capturing degradation trends that evolve across multiple sampling windows—a capability absent in purely convolutional architectures. Bidirectional LSTMs further improved temporal representation by processing sequences in both forward and backward directions.

Zhao et al. [3] introduced end-to-end CNN-LSTM architectures for fault diagnostics, demonstrating that the combination of convolutional feature extraction and LSTM temporal modeling outperformed either architecture alone on the CWRU bearing dataset. Zhang et al. [1] extended this with attention mechanisms, enabling the model to focus on the most diagnostic segments of each input window. Pham et al. [4] added bidirectional LSTM layers and attention, reporting 96.2% accuracy on a multi-fault rotating machinery dataset. Gupta and Sharma [5] specifically evaluated CNN-BiLSTM integration for motor fault diagnosis in industrial simulation environments, reporting strong performance across multiple fault severity levels.

D. Transformer and Attention-Based Methods

Vaswani et al. [8] introduced the Transformer architecture based entirely on self-attention mechanisms, which has subsequently been applied to time-series anomaly detection with notable success. Transformer-based models for time-series classification leverage multi-head self-attention to model long-range dependencies without the sequential computation bottleneck of RNNs. However, Transformer architectures typically require substantially larger training datasets to realize their performance advantages over hybrid CNN-LSTM models, making them less suitable for industrial deployments with limited labeled fault data.

E. SCADA and Industrial Monitoring Systems

Conventional SCADA systems provide real-time data acquisition, supervisory control, and threshold-based alarming. They excel at collecting high-volume sensor data, providing operator visibility via HMI screens, and executing automated control responses to extreme parameter deviations. However, they lack predictive intelligence: fault detection operates on static thresholds configured at commissioning, generating alarms only after parameters have already exceeded safe operating limits. Commercial vibration analyzers such as those offered by SKF, Brüel and Kjær, and Emerson provide more sophisticated condition monitoring but at costs prohibitive for small and medium industrial enterprises and without integrated deep learning inference.

Canizo et al. [13] demonstrated real-time predictive maintenance for wind turbines using big data frameworks, confirming the practical feasibility of cloud-based continuous fault monitoring at scale. Li et al.

[15] proposed directed acyclic graph networks combining CNN and LSTM layers for remaining useful life prediction, extending the application of hybrid architectures beyond binary fault detection to continuous health index estimation. Lim and Zohren [10] provided a comprehensive survey of deep learning approaches for time-series forecasting, documenting the superior performance of hybrid recurrent-convolutional architectures across diverse application domains.

F. Identified Gaps and Positioning of This Work

A systematic review of the existing literature reveals three persistent gaps: (1) most deep learning fault detection publications operate in an offline evaluation mode without demonstrating real-time deployment capability; (2) hardware-software integration—bridging model inference results to immediate physical alert mechanisms on the factory floor—is rarely addressed; and (3) comprehensive end-to-end latency profiling from sensor reading to operator alert is largely absent. This work directly addresses all three gaps by delivering a fully integrated, validated, real-time predictive maintenance system with sub-100 ms end-to-end response.

III. PROPOSED SYSTEM ARCHITECTURE

The proposed system introduces a hybrid CNN-LSTM deep learning framework that continuously monitors multivariate sensor streams from industrial machinery and delivers fault notifications through both digital and physical channels within 87 ms of fault manifestation. The architecture is organized into three tightly integrated functional layers operating in a continuous processing pipeline, as described below.

A. Layer 1: Sensor Data Acquisition

The acquisition layer comprises six sensor modalities mounted at critical observation points on the monitored machine. ADXL345 tri-axial MEMS accelerometers ($\pm 16g$, 13-bit resolution, 3200 Hz maximum sampling rate) are mounted on the motor bearing housings to capture vibration signatures associated with bearing defects, rotor imbalance, and shaft misalignment. DS18B20 digital temperature sensors ($\pm 0.5^\circ C$ accuracy) monitor motor winding and bearing housing thermal behavior. MPX5700AP analog pressure transducers monitor hydraulic and pneumatic circuit integrity. ACS712 Hall-effect current sensors capture motor phase current signatures, which carry rich fault information including turn-to-turn short circuits and supply voltage asymmetry. An MQ-135 gas sensor detects abnormal fumes indicative of insulation breakdown or lubricant combustion.

All six sensor channels are sampled synchronously at 100 Hz by the NodeMCU ESP8266 microcontroller and transmitted to the Flask server over Wi-Fi using a lightweight JSON-encoded HTTP POST protocol. Hardware timestamps are embedded in each packet to enable precise time-domain synchronization of all channels at the server. The acquisition layer is designed for hot-swap sensor replacement; failed transducers are detected through watchdog monitoring of expected packet arrival intervals.

B. Layer 2: Deep Learning Inference Engine

The inference engine is implemented as a Python process running alongside the Flask web server. It maintains a circular buffer that accumulates incoming sensor packets and triggers inference upon

completion of each 512-sample window with 50% overlap. The trained hybrid CNN-LSTM model, converted to TensorFlow Lite format for optimized inference, processes normalized sensor windows and outputs a five-class probability vector. The predicted class label and maximum probability score are persisted to a SQLite time-series database and simultaneously broadcast to all connected dashboard clients via WebSocket.

A confidence-gated alert mechanism suppresses fault notifications when the maximum class probability falls below a tunable threshold (default 0.85), preventing borderline uncertain predictions from generating nuisance alarms. When a fault is confirmed, the inference engine dispatches an HTTP POST alert payload containing fault class identity, confidence score, sensor channel readings, and timestamp to the NodeMCU hardware endpoint.

C. Layer 3: Alert Delivery and Operator Interface

The alert delivery layer operates through two parallel channels. The visual channel delivers fault status to the Flask web dashboard, where color-coded health indicators (green: normal, amber: warning confidence 0.70–0.84, red: fault confidence ≥ 0.85), live multi-channel waveform plots rendered by Chart.js, and a scrollable timestamped fault event log provide comprehensive situational awareness to remote operators. The audible channel drives the NodeMCU buzzer with a severity-adaptive pattern encoded in the alert payload: single beep for warning, intermittent for fault, continuous for critical. A 16×2 LCD module connected to the NodeMCU displays local status without requiring a networked display at the machine.

D. System Integration and Data Flow

Data flows continuously from physical sensors through the acquisition layer into the inference engine and out through the alert delivery layer. The NodeMCU serves a dual role: as an acquisition gateway transmitting sensor data upstream to Flask, and as an alert actuator receiving fault notifications downstream from Flask. This bidirectional communication architecture minimizes hardware by reusing the same Wi-Fi-connected microcontroller for both roles. The Flask server acts as the central intelligence hub coordinating acquisition, inference, logging, dashboard delivery, and hardware alert dispatch.

IV. METHODOLOGY

A. Experimental Testbed

The experimental testbed comprised a 2.2 kW three-phase induction motor (1440 RPM, 50 Hz, IP55 enclosure) coupled to an adjustable mechanical load via a flexible shaft coupling. Six sensor modalities—vibration (ADXL345), temperature (DS18B20), pressure (MPX5700AP), current (ACS712 ×3 for three phases), and air quality (MQ-135)—were instrumented at manufacturer-recommended measurement points. The testbed was operated in a climate-controlled laboratory environment maintained at $22^{\circ}\text{C} \pm 2^{\circ}\text{C}$ ambient temperature.

Five operational conditions were established and validated by a certified mechanical engineer: (1) Normal operation at rated load; (2) Bearing fault—artificially induced by removing four bearing balls from the drive-end bearing; (3) Overheating—induced by blocking 60% of cooling airflow using a physical baffle; (4) Pressure anomaly—

created by introducing a controlled hydraulic circuit leakage valve; (5) Electrical fault—simulated by interrupting one phase supply using a controlled contactor at the motor control center.

B. Dataset Construction

Continuous sensor data was recorded for 30 minutes per fault class at steady-state load conditions, yielding 180,000 samples per class at 100 Hz sampling rate. Raw signals were segmented using a sliding window of 512 samples (5.12 seconds) with 50% overlap stride, producing 12,400 labeled windows distributed equally across the five classes (2,480 per class). Class balance was maintained throughout to prevent classifier bias. Z-score normalization was applied independently per sensor channel using statistics computed from the training partition to prevent data leakage.

An IQR-based outlier filter with multiplier 3.0 removed isolated transient noise spikes attributable to electrical interference, affecting fewer than 0.3% of total samples. Derived statistical features—RMS amplitude, peak-to-peak value, crest factor, kurtosis, and skewness—were computed per window per channel and appended to the raw time-series representation as auxiliary input features. A stratified 80/20 train-test split (9,920/2,480 samples) preserved class proportions in both partitions. Five-fold stratified cross-validation was used during hyperparameter optimization.

C. CNN-LSTM Model Architecture

The hybrid model architecture processes input tensors of shape (512, 6)—512 time steps across 6 sensor channels—through two stages of feature extraction followed by temporal modeling and classification. The first convolutional block applies 64 filters of kernel size 3 with ReLU activation, batch normalization, and max pooling (pool size 2, stride 2), reducing the sequence length from 512 to 256 and extracting low-level local temporal features. The second convolutional block applies 128 filters of kernel size 3 with the same normalization and pooling sequence, further reducing to 128 time steps and learning higher-level discriminative patterns.

The 128-step feature sequences are passed to two stacked bidirectional LSTM layers, each with 128 hidden units per direction (256 total per layer). Bidirectional processing enables the model to exploit both past and future context within each window, improving sensitivity to fault signatures that are characterized by a specific temporal shape rather than a point anomaly. Residual connections between LSTM layers improve gradient flow and training stability. A dropout layer (rate 0.4) applied after the second LSTM layer regularizes the network against overfitting on the 9,920-sample training set.

The classification head consists of a dense layer with 128 units, batch normalization, ReLU activation, and a final softmax output layer with 5 neurons. The total parameter count is approximately 847,000, well within the memory constraints of the deployment server. The model was implemented in TensorFlow 2.12 with Keras functional API and converted to TensorFlow Lite format using post-training dynamic-range quantization, achieving a 4× reduction in model file size with less than 0.3% accuracy degradation.

D. Training Procedure and Hyperparameter Optimization

The model was trained using the Adam optimizer with an initial learning rate of 0.001. A ReduceLROnPlateau callback reduced the learning rate by factor 0.5 if validation loss did not improve for 5 consecutive epochs, enabling fine-grained convergence in later training stages. Early stopping with patience 15 epochs and best-weight restoration prevented overfitting. Categorical cross-entropy was used as the loss function. Training was performed with batch size 32 over a maximum of 200 epochs, achieving convergence at approximately 68 epochs.

Hyperparameter optimization was conducted using a randomized search over the following ranges: number of CNN filters (32, 64, 128), kernel sizes (3, 5, 7), LSTM hidden units (64, 128, 256), dropout rates (0.2, 0.3, 0.4, 0.5), and learning rates ($1e-4$, $5e-4$, $1e-3$). The final configuration was selected based on five-fold cross-validation F1-score. Training on an NVIDIA RTX 3060 GPU (12 GB VRAM, CUDA 11.8) required approximately 55 minutes for the complete training run.

E. Real-Time Inference Pipeline

At runtime, the NodeMCU ESP8266 samples all six sensor channels at 100 Hz, assembles readings into JSON payloads of 10 samples per packet, and transmits to the Flask server via HTTP POST at 10 Hz. The server accumulates packets in a thread-safe circular buffer and triggers a separate inference thread upon completion of each 512-sample window. Per-channel Z-score normalization uses running statistics updated with exponential moving averages to account for slow sensor drift during extended operation.

The TensorFlow Lite interpreter is pre-loaded into memory at server startup, eliminating model loading overhead from the inference path. Inference executes synchronously within the inference thread, completing in an average of 12 ms per window. Results are persisted to SQLite and broadcast via Flask-SocketIO to all connected browsers within 14 ms of inference completion. Fault alerts trigger an asynchronous HTTP POST to the NodeMCU endpoint, which processes the JSON payload and activates the buzzer within 12 ms of receipt, contributing to the 87 ms total end-to-end latency.

V. SYSTEM MODULES

A. Hybrid CNN-LSTM Feature Extraction Module

The CNN-LSTM feature extraction module implements the core analytical intelligence of the system. The 1D convolutional layers process multi-channel sensor input as a 2D spatial structure, simultaneously learning cross-channel correlations and within-channel temporal patterns. Filter bank analysis with 64 first-layer filters and 128 second-layer filters provides a rich 128-dimensional representation per time step after the second pooling layer. Batch normalization after each convolution stabilizes activations and accelerates convergence by reducing internal covariate shift.

The bidirectional LSTM layers perform temporal integration over the 128-step sequence produced by the convolutional trunk, encoding the temporal evolution of fault-relevant features across the full 5.12-second analysis window. The forward LSTM pass captures causal temporal dependencies, while the backward pass encodes contextual patterns whose diagnostic significance becomes apparent in retrospect. This dual-direction modeling is particularly valuable for gradual fault

types such as thermal degradation, where the fault signature accumulates progressively across the analysis window.

B. Multivariate Sensor Data Processing Module

The data processing module performs acquisition-time and preprocessing-time operations. At acquisition time, it validates packet integrity via CRC checksums embedded in the JSON payloads, detects and logs dropped packets, and performs timestamp-based inter-channel synchronization to within ± 1 sample. At preprocessing time, it applies IQR-based spike filtering, per-channel exponential moving average normalization, sliding window segmentation, and auxiliary feature computation (RMS, peak-to-peak, crest factor, kurtosis, skewness per channel). The module exposes a thread-safe queue interface to the inference engine, decoupling acquisition and processing from model inference.

C. Flask Web Application and Dashboard Module

The Flask web application serves as the system's central intelligence coordinator and human-machine interface. The application is structured with three distinct components: a REST API layer handling sensor data ingestion from the NodeMCU and providing historical fault query endpoints; a WebSocket event handler implemented via Flask-SocketIO managing real-time push of sensor readings and prediction results to connected browsers; and a static file server delivering the HTML/JavaScript dashboard application.

The dashboard renders six live sensor waveform plots using Chart.js with a 30-second rolling display window, updating at 10 Hz. A central health status panel displays the current machine state as a color-coded indicator with the predicted fault class label and confidence percentage. A fault event log table below the waveform panel maintains the 100 most recent fault events with timestamps, fault class, confidence scores, and sensor readings at detection time. A historical trend chart shows the fault probability time series over the preceding 24 hours.

D. Hardware Controller Module

The NodeMCU ESP8266 hardware controller firmware implements a multi-threaded event loop using the Arduino framework's cooperative scheduling model. The main loop handles sensor sampling at 100 Hz, packet assembly, and HTTP POST transmission to Flask. A parallel polling task queries the Flask fault notification endpoint at 100 ms intervals, parses incoming JSON payloads, and executes the appropriate actuator response. The buzzer driver implements three alert patterns: two short beeps for warning-class notifications (confidence 0.70–0.84), rapid intermittent beeping for fault-class notifications (confidence 0.85–0.94), and continuous sounding for critical-class notifications (confidence ≥ 0.95).

The LCD display module shows the current machine status on line 1 and the most recent fault class on line 2, refreshing at 1 Hz. EEPROM logging persists the timestamp, fault class, and confidence of each fault event in a circular buffer of 64 entries, preserving diagnostic history through power cycles and network outages. A hardware watchdog timer with a 30-second timeout automatically resets the microcontroller if any software thread stalls, ensuring continuous operation in unattended deployment scenarios.

E. Database and Logging Module

All fault events and sensor readings are persisted to a SQLite database structured around two primary tables: sensor_readings (timestamp, channel_id, value, window_id) and fault_events (timestamp, fault_class, confidence, window_id, sensor_snapshot_json). SQLite was selected for its zero-configuration deployment model and adequate performance at the 100 Hz sensor acquisition rate. Database writes are batched in transactions of 100 records to minimize write amplification and I/O overhead. A background archiving process compresses records older than 7 days into gzip-compressed CSV archives, maintaining the active database below 500 MB for indefinite operation.

VI. HARDWARE AND SOFTWARE REQUIREMENTS

A. Hardware Components

Table II presents all hardware components with their specifications and functional roles. The platform was selected to maximize capability within a strict cost budget suitable for small and medium industrial enterprises. Total bill-of-materials cost for a single monitoring node is approximately USD 45, enabling economical deployment across large fleets of machines.

TABLE II
Hardware Components, Specifications, and Functional Roles

Component	Specification	Func
Node MCU ESP8266	80 MHz Tensilica L106, 4MB Flash, Wi-Fi 802.11 b/g/n	Primary embedded Wi-Fi gateway
ADXL345	±16g, 3-axis, 13-bit, 3200 Hz max, I2C/SPI	Vibration & mechanical oscillation capture
DS18B20	-55°C to +125°C, ±0.5°C, 1-Wire protocol	Thermal behavior monitoring
MPX5700AP	15-700 kPa, 6.4 mV/kPa sensitivity, analog	Hydraulic/pneumatic sensing
MQ-135 Gas Sensor	10-300 ppm, analog output, 5V supply	Ambient air quality detection
ACS712 (20A)	±20A, 100 mV/A sensitivity, 5V supply	Motor current signature monitoring
Buzzer Module	85 dB SPL, active piezo, 5V DC	Audible fault alarm
16x2 LCD Module	HD44780 controller, I2C backpack, 5V	Local status display
Regulated PSU	5V/3.3V rails, 2A peak output	Stable multi-rail power
3-Phase Motor	2.2 kW, 1440 RPM, IP55, 50 Hz	Primary test asset

B. Software Stack

Table III presents the complete software stack. The entire system is implemented using open-source tools to maximize reproducibility and minimize licensing costs. All components run on standard Linux (Ubuntu 20.04) or Windows 10 host systems without specialized operating system requirements.

TABLE III
Software Stack with Version and Role

Tool/Library	Version	Role
Python	3.10	Primary programming language
TensorFlow / Keras	2.12	Deep learning model training & inference
TensorFlow Lite	2.12	Optimized edge model deployment
Flask	2.3	HTTP server & REST API routing
Flask-SocketIO	5.3	WebSocket real-time push to dashboard
NumPy	1.24	Numerical array operations
Pandas	2.0	Data preprocessing & time-series management
scikit-learn	1.3	Evaluation metrics, stratified splits
Matplotlib / Seaborn	3.7 / 0.12	Training visualization & confusion matrices
Chart.js	4.x	Browser-based live sensor waveform plot
SQLite	3.x	Time-series fault event logging
Arduino IDE	2.1	NodeMCU ESP8266 firmware development
ArduinoJson	6.x	JSON payload parsing on microcontroller
Git	2.x	Version control for all source code

tion

ed controller &

C. Network and Communication Architecture

The system employs a star topology Wi-Fi network with the Flask server acting as the central hub. The NodeMCU connects to the server's Wi-Fi access point using WPA2 security. Sensor data is transmitted via HTTP POST at 10 Hz (average 4.8 KB/s upstream bandwidth). Fault alert polling uses HTTP GET at 10 Hz (average 0.2 KB/s downstream bandwidth). Dashboard clients connect via WebSocket over HTTP/1.1 with TLS encryption when accessed remotely. The total network bandwidth requirement is below 100 KB/s per monitoring node, compatible with standard industrial Wi-Fi infrastructure.

II. EXPERIMENTAL RESULTS AND PERFORMANCE EVALUATION

A. Per-Class Classification Performance

The hybrid CNN-LSTM model was evaluated on the held-out test set of 200 samples using per-class precision, recall, F1-score, and support. Table IV presents detailed per-class metrics. The model achieved the highest performance on Normal Operation (F1: 98.5%) and Bearing Failure (F1: 98.1%) classes, both characterized by highly distinct sensor signatures. Overheating and Electrical Fault show marginally lower performance (F1: 95.9% and 96.1% respectively) due to overlapping thermal signatures in the temperature sensor channel, which is the primary discriminative channel for both classes.

TABLE IV
Per-Class Classification Metrics on Test Set (n=2,480)

Fault Class	Prec. (%)	Rec. (%)	F1 (%)	Support
Normal Operation	98.2	98.9	98.5	496
Bearing Failure	97.9	98.4	98.1	496

Overheating	96.1	95.8	95.9	
Pressure Anomaly	97.3	97.0	97.1	
Electrical Fault	96.5	95.7	96.1	
Macro Average	97.2	97.2	97.1	

B. Comparative Analysis Against Baselines

Table I presents comprehensive comparative performance against eight baseline methods. The SVM+FFT baseline [2] achieved 87.5% accuracy with a 9.3% false positive rate. kNN with wavelet features achieved 85.2% accuracy. Standalone CNN achieved 91.2%; standalone LSTM achieved 93.7%. CNN with attention mechanism achieved 94.8%; Bidirectional LSTM alone achieved 95.1%. The conventional threshold-based SCADA system achieved only 78.4% accuracy with a 14.8% false positive rate—demonstrating the substantial performance gap between conventional industrial practice and deep learning approaches. The proposed CNN-LSTM model achieved 97.4% accuracy with 1.8% false positive rate, the highest performance across all evaluated methods.

TABLE I

Comparative Performance of Fault Detection Methods on Test Set

Method	Acc. (%)	Prec. (%)	Rec. (%)
SVM+FFT [2]	87.5	85.3	86.1
kNN+Wavelet	85.2	83.1	84.0
CNN only	91.2	90.1	90.8
LSTM only	93.7	92.4	93.1
Random Forest	89.1	87.9	88.5
CNN-Attention	94.8	93.9	94.5
Bi-LSTM	95.1	94.2	94.9
SCADA Threshold	78.4	75.0	77.2
Proposed CNN-LSTM	97.4	96.8	97.1

C. Ablation Study

An ablation study was conducted to quantify the contribution of individual architectural components. Removing the bidirectional processing from LSTM layers (replacing with unidirectional) reduced accuracy from 97.4% to 95.8%. Removing the second convolutional block reduced accuracy to 94.3%. Removing all convolutional blocks (LSTM-only) reduced accuracy to 93.7%. Removing auxiliary statistical features reduced accuracy by 0.6%. Removing batch normalization reduced accuracy by 1.1% and increased training time by approximately 40%. These results confirm that each architectural element contributes meaningfully to the final system performance.

D. End-to-End Latency Profiling

Table V presents detailed end-to-end latency measurements across all pipeline stages, measured over 1,000 consecutive inference cycles during a 4-hour continuous operation session. The dominant latency contributor is Wi-Fi transmission (18 ms average), followed by WebSocket dashboard push (14 ms) and HTTP POST to NodeMCU (16 ms). Total end-to-end latency from sensor sample to operator alert

averages 87 ms, well within the operational safety requirements of rotating machinery where fault progression timescales are measured in seconds to minutes.

TABLE V

End-to-End Pipeline Latency Profiling (n=1,000 measurements)

Pipeline Stage	Avg (ms)	P95 (ms)	P99 (ms)
Sensor → Server (Wi-Fi)	18	24	31
Buffer Fill (512 samples)	5	6	8
Z-score Normalization	2	3	4
TFLite Inference	12	15	18
DB Logging (SQLite)	8	11	14
WebSocket Push	14	19	25
HTTP POST to NodeMCU	16	21	28
Buzzer Activation	12	16	21
Total End-to-End	87	115	144

E. 72-Hour Reliability Test

A 72-hour continuous operation reliability test was conducted with fault events injected at random intervals (minimum 45-minute separation) by the supervising engineer. The system detected all 47 injected fault events, achieving a false negative rate of 0.0%. False alarm count was 8 over 72 hours of continuous operation (false positive rate 1.8% by event count). System uptime was 100% with the NodeMCU watchdog timer triggering one automatic reset at hour 51 due to a temporary Wi-Fi disconnection. CPU utilization on the Flask server remained below 18% throughout, and memory usage was stable at 312 MB.

F. Model Robustness Under Load Variation

To evaluate model robustness under varying operational loads, the testbed was operated at 25%, 50%, 75%, and 100% of rated load while injecting each fault class. Classification accuracy degraded modestly from 97.4% at rated load to 95.2% at 25% load, primarily due to reduced signal-to-noise ratio in the vibration and current channels at light load conditions. This performance was still substantially superior to all evaluated baselines, confirming the model's generalization capability across realistic load profiles.

G. Cross-Validation Results

Five-fold stratified cross-validation on the complete 12,400-sample dataset yielded a mean accuracy of 97.1% ± 0.4% (mean ± standard deviation across folds), confirming stable generalization with low variance across data partitions. The lowest fold accuracy was 96.5%, and the highest was 97.8%. These results confirm that the reported test set accuracy is representative and not an artifact of a particularly favorable data split.

VIII. DISCUSSION

A. Interpretation of Results

The experimental results establish that the hybrid CNN-LSTM architecture provides substantially superior fault detection performance compared to all evaluated baselines across both accuracy

and false positive rate dimensions. The 97.4% accuracy achieved on a balanced five-class dataset with realistic industrial noise represents a significant improvement over the 78.4% accuracy and 14.8% false positive rate of conventional threshold-based SCADA—the type of system currently deployed in most industrial facilities. The practical implication is a potential 8.5-fold reduction in false alarms compared to existing practice, directly addressing the alarm fatigue problem that reduces operator responsiveness in industrial control rooms.

The architectural ablation study confirms that both the convolutional and recurrent components make necessary contributions: the CNN blocks provide essential local feature extraction that the LSTM alone cannot replicate, while the LSTM temporal modeling provides fault persistence discrimination that distinguishes genuine degradation from transient disturbances. The bidirectional LSTM configuration provides a 1.6 percentage point advantage over its unidirectional equivalent, confirming that backward temporal context contributes meaningfully to fault characterization.

B. Operator Interface Evaluation

User acceptance testing was conducted with four plant maintenance technicians and two shift supervisors at the test facility. Participants were asked to monitor the system dashboard during controlled fault injection sessions without prior knowledge of injection timing or fault type. All six participants successfully identified fault events from the dashboard color-coded indicators and waveform anomalies within an average of 12 seconds of system detection. Participants rated the dashboard interface 4.3/5.0 on usability and 4.5/5.0 on information clarity.

Qualitative feedback highlighted the confidence score display as particularly valuable: operators could distinguish between high-confidence critical faults requiring immediate response and lower-confidence borderline events warranting scheduled inspection. The fault event log with historical confidence trends enabled operators to track degradation progression over time, supporting planning of maintenance interventions at the most economical juncture.

C. Identified Limitations

Several limitations were identified through extended validation and operator feedback. First, sensor drift: over the 72-hour reliability test, gradual drift in the MQ-135 ambient gas sensor baseline was observed, causing a marginal increase in false alarm rate from 1.4% in the first 24 hours to 2.3% in hours 48–72. The current Z-score normalization using exponential moving average statistics partially compensates for drift but does not fully correct for large-scale sensor aging across months of operation.

Second, Wi-Fi reliability: the system assumes stable Wi-Fi connectivity between the NodeMCU and the Flask server. During the 72-hour test, one connectivity interruption lasting 23 seconds was observed at hour 51, triggering the watchdog reset. In industrial facilities with extensive metallic infrastructure, RF signal attenuation and interference from variable-frequency motor drives could reduce Wi-Fi reliability below acceptable levels, necessitating a wired Ethernet fallback.

Third, fault class coverage: the current five-class taxonomy covers the most common rotating machinery faults but does not include compound faults (simultaneous occurrence of multiple fault types), partial discharge in windings, cavitation in pumps, or gearbox tooth damage. Extending the fault class taxonomy requires collection of additional labeled data for each new class, which is time-consuming and operationally disruptive to production facilities.

Fourth, scalability: the current single-node architecture supports one machine per Flask server instance. Monitoring a fleet of 50 machines would require either a distributed architecture with multiple server instances or migration to a cloud-based inference platform with load balancing. The SQLite database also presents throughput limitations above approximately 10 concurrent sensor nodes.

D. Comparison with Commercial Solutions

Commercial predictive maintenance platforms such as SKF Enlight AI, Emerson AMS Device Manager, and Honeywell Forge offer sophisticated machine learning-based condition monitoring capabilities but at subscription costs of \$500–\$5,000 per machine per year. The proposed system's bill-of-materials cost of approximately USD 45 per node with zero ongoing licensing cost represents a cost reduction factor of 100–1,000× compared to commercial alternatives. For small and medium industrial enterprises operating fleets of 10–100 machines, this cost difference is transformative for the business case of predictive maintenance adoption.

E. Energy and Environmental Considerations

The NodeMCU ESP8266 consumes approximately 170 mA at 3.3V during active Wi-Fi transmission, corresponding to 0.56 W per node. The Flask server requires approximately 45 W of CPU power during continuous operation. For a 50-machine installation, total system power consumption would be approximately 73 W—a negligible fraction of the power consumed by the monitored machinery and readily powered from standard facility supply circuits. The open-source software stack eliminates electronic waste from proprietary hardware replacements at software end-of-life.

IX. FUTURE RESEARCH DIRECTIONS

A. Online Learning and Sensor Drift Adaptation

Maintaining model accuracy under gradual sensor drift and evolving machine condition is the most critical near-term research challenge. Online learning strategies—including continual learning with elastic weight consolidation and prototype-based class updating—will be investigated to enable the model to update its decision boundaries incrementally from streaming unlabeled data using pseudo-label confidence gating. Adaptive normalization layers that learn sensor-specific drift correction parameters from historical statistics will be integrated into the preprocessing pipeline.

B. Federated Learning for Multi-Site Deployment

Industrial operators are often reluctant to share proprietary operational data with external parties for model training. Federated learning frameworks enable collaborative model improvement across multiple site deployments without centralizing sensitive operational data. Each site trains a local model update on its own data, and only the gradient update (not the raw data) is shared with a central aggregation server. This approach will be implemented using the TensorFlow

Federated framework to enable cross-facility model improvement while preserving data sovereignty.

C. Transformer-Based Temporal Modeling

Transformer architectures with multi-head self-attention have demonstrated state-of-the-art performance on time-series classification benchmarks and will be evaluated as replacements for the LSTM temporal modeling component. The Temporal Fusion Transformer (TFT) architecture is of particular interest given its explicit handling of static and time-varying covariates, which maps naturally to the combination of fixed machine parameters and dynamic sensor readings in the proposed system. Attention weight visualization will also be explored to improve model interpretability for maintenance engineers.

D. Edge TPU Deployment

Deployment of the inference model on Google Coral Edge TPU accelerators will be investigated to reduce inference latency below 5 ms, enabling the system to respond to transient faults that manifest over sub-second timescales. Edge TPU deployment also reduces dependence on a networked server, enabling autonomous fault detection at the sensor node level—critical for deployments in remote locations or facilities with unreliable network connectivity.

E. Remaining Useful Life Prediction

Extending the system from binary fault detection to continuous remaining useful life (RUL) prediction represents a high-value research direction. RUL prediction enables operators to quantify the time remaining before a fault will reach a critical severity level, supporting optimal maintenance scheduling. Li et al. [15] demonstrated the feasibility of CNN-LSTM approaches for RUL prediction; this architecture will be extended with a regression output head and validated against the NASA C-MAPSS turbofan engine degradation dataset as a benchmark before deployment on the industrial motor testbed.

F. Explainability and Model Interpretability

Industrial operators and process engineers require not only accurate fault predictions but also actionable explanations: which sensor channels and which time intervals within a detection window contributed most to the fault classification decision. Gradient-weighted Class Activation Mapping (Grad-CAM) adapted for 1D convolutional layers will be implemented to generate per-sensor channel saliency maps for each detection event. SHAP (SHapley Additive exPlanations) will be applied to the final dense layer to quantify feature importance contributions from both raw and derived sensor features.

X. CONCLUSION

This paper has presented a comprehensive intelligent industrial machine fault detection system integrating a hybrid bidirectional CNN-LSTM deep learning architecture, a real-time Flask web monitoring dashboard with WebSocket communication, a six-channel multivariate sensor acquisition layer, and a NodeMCU ESP8266 hardware alert controller. The system architecture addresses the complete workflow from raw sensor data acquisition through deep learning inference to immediate physical operator notification, closing the loop between digital model intelligence and real factory-floor response.

Exhaustive experimental validation on a three-phase induction motor testbed under five representative fault conditions demonstrated 97.4% overall classification accuracy, precision and recall exceeding 96% across all five fault classes, a false positive rate of 1.8%, and an end-to-end system latency of 87 ms. All metrics were significantly superior to eight evaluated baselines including threshold-based SCADA, classical machine learning methods, and standalone deep learning architectures. Architectural ablation confirmed the essential contribution of each model component. A 72-hour continuous reliability test with 47 injected fault events demonstrated 0% false negative rate and 100% system uptime.

End-to-end latency profiling across all pipeline stages confirmed production suitability for the timescales of rotating machinery fault progression. The complete hardware bill-of-materials cost of approximately USD 45 per monitoring node positions the system as a cost-effective alternative to commercial predictive maintenance platforms, enabling adoption by small and medium industrial enterprises for whom commercial solutions are economically inaccessible.

Future research directions include online learning for sensor drift adaptation, federated learning for multi-site deployment, Transformer-based temporal modeling, edge TPU acceleration, remaining useful life prediction, and SHAP-based model interpretability. The proposed system provides a validated, deployable foundation for fully autonomous, AI-driven predictive maintenance in Industry 4.0 manufacturing environments, contributing meaningfully to the broader goals of reduced industrial downtime, improved worker safety, and more efficient utilization of industrial assets.

ACKNOWLEDGMENT

The authors gratefully acknowledge the guidance of Mentor Priyadharshan Sir and the support of the Department of Artificial Intelligence and Data Science, Nehru Institute of Engineering and Technology, Thirumalayampalayam, Tamil Nadu 641105, India. The authors thank the laboratory technicians for their assistance in setting up and instrumenting the industrial motor testbed, and the plant operators who participated in the user acceptance testing sessions.

REFERENCES

- [1] X. Zhang, Y. Li, and M. Chen, "Machine Fault Detection Using a Hybrid CNN-LSTM Attention-Based Model," *IEEE Trans. Ind. Electron.*, vol. 70, no. 4, pp. 3812-3821, Apr. 2023.
- [2] H. Wang, J. Liu, and R. Patel, "Deep Learning Based Fault Detection and Diagnosis Method for Power Systems," *IEEE Access*, vol. 13, pp. 10452-10463, Jan. 2025.
- [3] S. Zhao, T. Wu, and F. Dong, "End-to-End CNN + LSTM Deep Learning Approach for Fault Diagnostics," *Mech. Syst. Signal Process.*, vol. 120, pp. 514-527, Apr. 2019.
- [4] K. Pham, A. Nguyen, and D. Le, "Fault Detection with a 1D CNN + Bi-LSTM Attention Hybrid Model," *Sensors*, vol. 25, no. 3, p. 874, Feb. 2025.
- [5] R. Gupta and P. Sharma, "Integration and Simulation of a CNN-Bi-LSTM Model for Motor Fault Diagnosis," *J. Manuf. Syst.*, vol. 74, pp. 289-301, Mar. 2025.
- [6] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436-444, May 2015.
- [7] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735-1780, Nov. 1997.

- [8] A. Vaswani et al., "Attention is all you need," in Proc. NeurIPS, 2017, pp. 5998-6008.
- [9] Z. Chen, K. Gryllias, and W. Li, "Intelligent fault diagnosis for rotary machinery using transferable convolutional neural network," IEEE Trans. Ind. Inform., vol. 16, no. 1, pp. 339-349, Jan. 2020.
- [10] B. Lim and S. Zohren, "Time-series forecasting with deep learning: A survey," Philos. Trans. R. Soc. A, vol. 379, no. 2194, Apr. 2021.
- [11] Y. Lei, B. Yang, X. Jiang, F. Jia, N. Li, and A. K. Nandi, "Applications of machine learning to machine fault diagnosis: A review and roadmap," Mech. Syst. Signal Process., vol. 138, p. 106587, Apr. 2020.
- [12] W. Zhang, G. Peng, C. Li, Y. Chen, and Z. Zhang, "A new deep learning model for fault diagnosis with good anti-noise and domain adaptation ability on raw vibration signals," Sensors, vol. 17, no. 2, p. 425, Feb. 2017.
- [13] A. Canizo, E. Onieva, A. Conde, S. Charramendieta, and S. Trujillo, "Real-time predictive maintenance for wind turbines using BigData frameworks," in Proc. IEEE ICPHM, 2017, pp. 60-67.
- [14] R. Zhao, R. Yan, Z. Chen, K. Mao, P. Wang, and R. X. Gao, "Deep learning and its applications to machine health monitoring," Mech. Syst. Signal Process., vol. 115, pp. 213-237, Jan. 2019.
- [15] J. Li, X. Li, and D. He, "A directed acyclic graph network combined with CNN and LSTM for remaining useful life prediction," IEEE Access, vol. 7, pp. 75365-75379, 2019.
- [16] N. Christiansen, A. Graening, F. Rotgeri, M. Lipsmeier, and J. Beyerer, "Performance evaluation of unsupervised anomaly detection methods for industrial applications," in Proc. ECML PKDD Workshop, 2021.
- [17] P. Nectoux et al., "PRONOSTIA: An experimental platform for bearings accelerated degradation tests," in Proc. IEEE ICPHM, 2012.
- [18] K. A. Loparo, "Bearings Vibration Data Set," Case Western Reserve University, Cleveland, OH, USA, [Online]. Available: <https://engineering.case.edu/bearingdatacenter>.
- [19] O. Janssens et al., "Convolutional neural network based fault detection for rotating machinery," J. Sound Vib., vol. 377, pp. 331-345, Sep. 2016.
- [20] A. Gretton, K. Borgwardt, M. Rasch, B. Scholkopf, and A. Smola, "A kernel two-sample test," J. Mach. Learn. Res., vol. 13, pp. 723-773, Mar. 2012.
- [21] P. Rodriguez, J. Mendez-Rebolledo, and R. Saenz-Aguirre, "Fault detection in induction motors based on MCSA and stacked autoencoders," IEEE Trans. Ind. Electron., vol. 68, no. 3, pp. 2100-2110, Mar. 2021.
- [22] M. Gan, C. Wang, and C. A. Zhu, "Construction of hierarchical diagnosis network based on deep learning and its application in the fault pattern recognition of rolling element bearings," Mech. Syst. Signal Process., vol. 72-73, pp. 92-104, May 2016.
- [23] H. Shao, H. Jiang, H. Zhang, and T. Liang, "Electric locomotive bearing fault diagnosis using a novel convolutional deep belief network," IEEE Trans. Ind. Electron., vol. 65, no. 3, pp. 2727-2736, Mar. 2018.
- [24] F. Jia, Y. Lei, J. Lin, X. Zhou, and N. Lu, "Deep neural networks: A promising tool for fault characteristic mining and intelligent diagnosis of rotating machinery with massive data," Mech. Syst. Signal Process., vol. 72-73, pp. 303-315, May 2016.
- [25] X. Li, W. Zhang, Q. Ding, and J.-Q. Sun, "Intelligent rotating machinery fault diagnosis based on deep learning using data augmentation," J. Intell. Manuf., vol. 31, pp. 433-452, Feb. 2020.