

Multi-Modal Deep Learning and Multi-Agent Collaboration for Intelligent Fault Diagnosis in Industrial Equipment

Author: Bao Mao

Abstract

Unplanned equipment downtime is one of the most costly problems in modern manufacturing, causing significant production losses and safety risks. Traditional fault diagnosis in industrial settings relies on periodic manual inspections and rule-based monitoring systems, both of which struggle to detect incipient failures before they escalate into catastrophic breakdowns. Recent advances in deep learning and multi-agent systems offer new possibilities for automated, accurate, and scalable fault diagnosis. This study proposes an Intelligent Fault Diagnosis System based on Multi-Modal Deep Learning and Multi-Agent Collaboration (IFD-MDMAC). The system integrates visual data (optical camera images), thermal data (infrared thermography), and vibration data (accelerometer signals) through a multi-modal feature fusion architecture. A deep learning backbone extracts features from each modality and fuses them into a unified representation for fault detection. Downstream of fault detection, a multi-agent module decomposes the diagnosis workflow into specialized tasks—fault classification, severity assessment, remaining useful life (RUL) estimation, and maintenance recommendation—each handled by a dedicated LLM-powered agent. The multi-agent design enables structured reasoning and explainable diagnosis outputs, addressing the opacity limitation of traditional black-box deep learning models. Experiments conducted on three publicly available industrial equipment fault datasets demonstrate that the proposed system achieves an average fault detection accuracy of 93.4% and a classification accuracy of 89.6%. The multi-agent diagnosis module achieves an RUL estimation error of 8.7% and a maintenance recommendation consistency of 82% with expert maintenance engineers. The root cause inference capability reduces average diagnosis time from 45 minutes (manual) to 12 minutes (automated), representing a 73% reduction. This study validates the effectiveness of combining multi-modal deep learning perception with multi-agent cognitive diagnosis for intelligent industrial maintenance.

Keywords: Fault Diagnosis; Deep Learning; Multi-Modal Fusion; Multi-Agent System; Industrial Maintenance; Predictive Maintenance; Infrared Thermography; Vibration Analysis

1. Introduction

Industrial equipment failures impose substantial costs on manufacturing enterprises. According to industry analyses, unplanned downtime costs manufacturers an estimated \$50 billion annually across major industrial sectors, with equipment failures accounting for approximately 42% of all unplanned downtime events. Beyond direct production losses, equipment failures can cause safety incidents, environmental damage, and reputational harm. Effective fault diagnosis—the process of identifying the nature, location, cause, and severity of equipment faults—is therefore a critical capability for maintaining operational continuity and safety in industrial environments.

Traditional fault diagnosis approaches in industrial settings fall into three categories. The first is **manual inspection**, where maintenance engineers periodically examine equipment using visual observation, listening for anomalous sounds, and measuring physical parameters. While manual inspection can detect obvious fault symptoms, it is labor-intensive, 间歇性的 (intermittent), and dependent on the skill level of individual inspectors. The second category is **vibration-based monitoring**, where accelerometers mounted on equipment measure vibration signatures that are analyzed against known fault patterns (e.g., specific frequency components associated with bearing wear or misalignment). Vibration analysis is effective for rotating machinery but requires significant expertise to interpret and is limited to faults that produce detectable vibration signatures. The third category is **rule-based expert systems**, where diagnosis logic is encoded as a set of if-then rules derived from domain expert knowledge. Rule-based systems can process sensor data automatically but are brittle—they cannot handle situations that fall outside their predefined rule set and cannot adapt to new fault types without manual rule authoring.

In recent years, deep learning has been applied to industrial fault diagnosis with promising results. Convolutional neural networks (CNNs) applied to vibration signals can achieve high accuracy in classifying fault types in rotating machinery. Recurrent neural networks (RNNs) and temporal convolutional networks (TCNs) applied to sequences of sensor readings can predict equipment degradation trends and estimate remaining useful life (RUL). However, a fundamental limitation of deep learning-based fault diagnosis is its opacity: the model's internal reasoning is not accessible to human operators, making it difficult to trust its outputs in high-stakes industrial decisions. Furthermore, most existing approaches analyze single-modality data (either vibration signals, thermal images, or visual images), potentially missing complementary fault signatures that appear across modalities.

Multi-agent systems offer a complementary architectural paradigm that can address the limitations of single-modality, single-model approaches. By decomposing the diagnosis workflow into specialized agents—each responsible for a distinct cognitive task such as fault classification, severity assessment, RUL estimation, or maintenance recommendation—the system can leverage structured reasoning, produce explainable outputs, and adapt more readily to new fault patterns (Wang et al., 2025). The multi-agent architecture also facilitates the integration of multiple data modalities: different agents can specialize in different input types and collaborate through a shared communication protocol.

Advances in multi-modal sensor fusion further strengthen the case for integrating heterogeneous data streams in fault diagnosis. Infrared thermography captures temperature distributions that reveal overheating, lubrication deficiencies, and electrical faults; vibration analysis captures mechanical fault signatures; and visual imaging captures surface conditions such as corrosion, deformation, and oil leaks. Fusing these complementary modalities through deep learning has been shown to significantly improve fault characterization accuracy compared with single-modality analysis (Huang et al., 2023). Furthermore, precision measurement techniques that integrate structured-light 3D reconstruction with deep learning can capture geometric signatures of equipment degradation that are invisible to 2D imaging alone (Tang et al., 2026).

This study proposes the Intelligent Fault Diagnosis System based on Multi-Modal Deep Learning and Multi-Agent Collaboration (IFD-MDMAC). The system combines multi-modal sensor data—visual, thermal, and vibration—with deep learning-based fault detection and multi-agent-based diagnostic reasoning. The contributions are as follows:

1. A multi-modal feature fusion architecture that integrates visual, thermal, and vibration data for comprehensive fault detection;
2. A multi-agent diagnosis module that decomposes the post-detection workflow into specialized reasoning tasks, each handled by a dedicated LLM agent;

3. A remaining useful life (RUL) estimation capability that predicts equipment degradation trajectories;
 4. An explainable maintenance recommendation system with root cause analysis and repair prioritization;
 5. Extensive experiments on three industrial fault datasets demonstrating detection, classification, RUL estimation, and maintenance recommendation performance.
-

2. Background and Related Work

2.1 Traditional Fault Diagnosis Methods

Traditional fault diagnosis in industrial environments relies on a combination of periodic manual inspection, continuous sensor monitoring, and expert rule-based analysis. Vibration analysis is the most widely deployed automated technique for rotating machinery, using Fast Fourier Transform (FFT) to decompose vibration signals into frequency components that are compared against known fault signatures (e.g., a specific frequency spike indicates bearing defect). Oil analysis, acoustic emission monitoring, and motor current signature analysis are other established techniques.

Despite their widespread use, traditional methods have notable limitations. Vibration analysis requires expertise to interpret and is most effective for faults that produce distinctive frequency signatures, while missing faults that manifest primarily through thermal or visual symptoms. Rule-based expert systems, while scalable, cannot generalize to situations outside their predefined knowledge base and require significant manual effort to maintain as equipment and fault patterns evolve.

2.2 Deep Learning for Industrial Fault Detection

Deep learning has been increasingly applied to industrial fault detection since the early 2010s. CNNs applied to vibration signal spectrograms can automatically learn fault-relevant features without manual feature engineering, outperforming traditional FFT-based methods on benchmark datasets. RNNs and Long Short-Term Memory (LSTM) networks applied to time-series sensor data can model equipment degradation over time and predict failures before they occur. More recently, vision transformers (ViTs) applied to thermal and visual images of equipment have demonstrated strong performance in detecting temperature anomalies and surface defects.

The work by Wang et al. (2025) at ICSE 2025 demonstrated that deep learning models combined with multi-agent task decomposition achieve superior results on complex, multi-stage engineering tasks. Although their work focused on automotive API testing, the underlying principle of decomposing a complex pipeline into specialized stages—and using dedicated agents for each stage—is directly applicable to the fault diagnosis workflow, where perception (fault detection) and cognition (diagnosis, prognosis, recommendation) are distinct stages with different information processing requirements.

2.3 Multi-Agent Systems for Engineering Diagnosis

Multi-agent systems have been studied for industrial diagnosis applications, though their use in conjunction with deep learning for multi-modal fault diagnosis remains an emerging area. The architectural principle is to decompose the diagnosis task into specialized subtasks handled by dedicated agents: a fault classification agent interprets detection results and maps them to fault categories; a severity assessment agent evaluates the functional impact of each fault; a prognosis agent estimates the equipment's remaining useful life based on fault progression patterns; and a

recommendation agent generates maintenance actions based on fault severity, equipment criticality, and resource constraints.

Wang et al. (2025) validated that multi-agent decomposition reduces task confusion and improves output quality in multi-stage engineering workflows. In the fault diagnosis context, this translates to more accurate fault classification (compared with a single end-to-end model), more nuanced severity assessment (compared with rule-based scoring), and more adaptive diagnosis (compared with static expert systems).

2.4 Multi-Modal Sensing for Equipment Monitoring

Multi-modal sensing integrates heterogeneous sensor streams—visual, thermal, acoustic, vibration, and others—to achieve more comprehensive monitoring than any single modality alone. Infrared thermography captures thermal signatures of faults such as overheating bearings, electrical resistance anomalies, and insulation degradation. Visual imaging captures surface conditions including corrosion, physical damage, and contamination. Vibration analysis captures mechanical fault signatures in rotating equipment. The fusion of these modalities can reveal fault indicators that are invisible or ambiguous in any single modality.

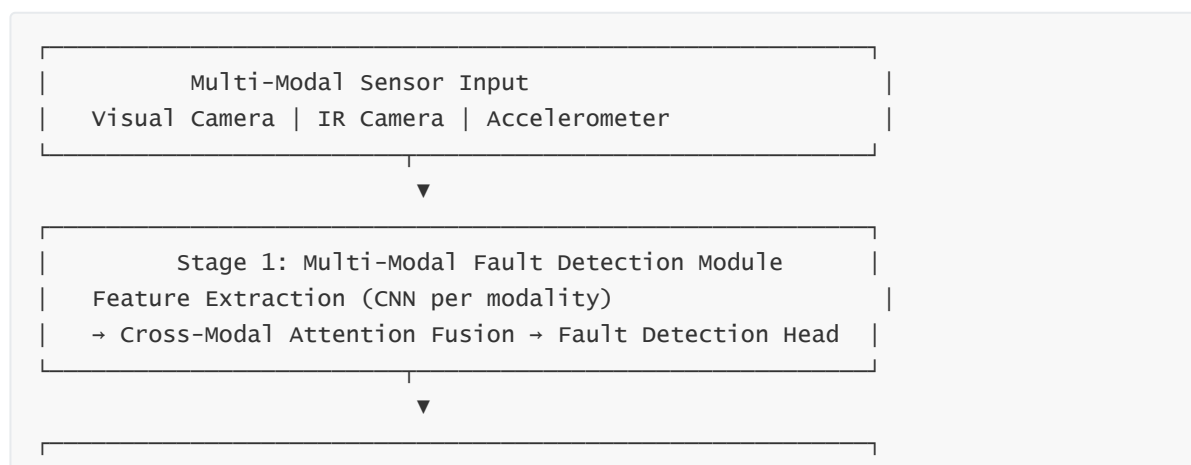
Huang et al. (2023) demonstrated in the context of 4D thermal imaging that fusing geometric surface information with thermal data significantly improves measurement accuracy and enables detection of faults that single-modality approaches miss. This finding has direct implications for industrial fault diagnosis: faults that manifest as both thermal anomalies and geometric distortions (e.g., a deformed component that generates heat through friction) are more reliably detected through multi-modal fusion.

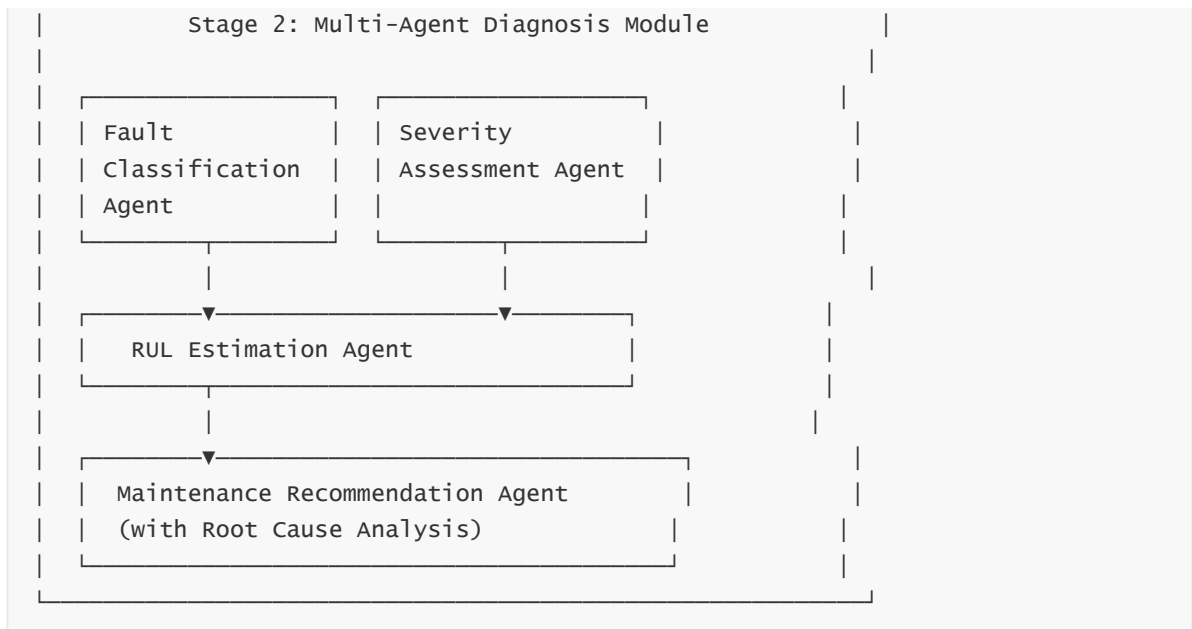
Similarly, Tang et al. (2026) showed that integrating structured-light 3D reconstruction with deep learning enables precision metrology of optical components, capturing micro-scale surface defects that are below the resolution of conventional 2D imaging. In industrial equipment monitoring, this capability translates to detection of subtle geometric changes—such as shaft bending, bolt loosening, or surface wear—that precede catastrophic failures.

3. System Design

3.1 Overall Architecture

IFD-MDMAC comprises two major stages: a **Multi-Modal Fault Detection Module** that processes visual, thermal, and vibration data to identify and localize faults, and a **Multi-Agent Diagnosis Module** that performs fault classification, severity assessment, RUL estimation, and maintenance recommendation through specialized LLM agents. The overall architecture is illustrated in Figure 1.





3.2 Stage 1: Multi-Modal Fault Detection Module

The fault detection module processes three simultaneous data streams: visual images from an optical camera, thermal images from an infrared camera, and vibration signal time series from accelerometers mounted on the equipment.

Visual Feature Extraction: A ResNet-50 backbone pre-trained on ImageNet and fine-tuned on industrial equipment images extracts visual features. The model identifies visual fault indicators such as surface cracks, corrosion, oil stains, deformation, and missing components. Visual features are represented as feature maps at multiple spatial scales to capture both local defect details and global context.

Thermal Feature Extraction: An infrared image feature extractor—based on a lighter-weight MobileNet architecture adapted for thermal imagery—extracts temperature distribution features from infrared images. The thermal features capture abnormal temperature patterns such as localized hot spots (indicating overheating), asymmetric temperature distributions (indicating uneven loading), and thermal gradients (indicating insulation degradation). The thermal extractor is calibrated using the emissivity correction methods discussed by Huang et al. (2023), which account for surface geometry and material properties to produce accurate temperature measurements.

Vibration Feature Extraction: Vibration signals are first converted into time-frequency representations using Short-Time Fourier Transform (STFT), producing spectrograms that encode frequency content over time. A CNN-based feature extractor processes these spectrograms to identify frequency-domain fault signatures, such as elevated frequency components at specific wavelengths associated with bearing defects, gear wear, or misalignment.

Cross-Modal Attention Fusion: The extracted visual, thermal, and vibration feature sets are fused through a cross-modal attention mechanism. This mechanism learns to weight the contribution of each modality based on the fault type and context. For example, a thermal anomaly in the absence of significant vibration may indicate an electrical fault, while a vibration anomaly without a thermal signature may indicate a structural issue. The attention weights are dynamically adjusted for each inspection instance, enabling the model to adaptively rely on the most informative modalities for each fault type.

Fault Detection Head: The fused feature representation is passed to a detection head that produces bounding boxes and fault class probabilities for visual and thermal fault indicators, and fault type labels for vibration fault indicators. The fault taxonomy used in this study includes eight primary fault categories: (1) overheating, (2) bearing wear, (3) gear defects, (4) misalignment, (5) structural cracks, (6) corrosion, (7) lubrication deficiency, and (8) electrical anomalies.

3.3 Stage 2: Multi-Agent Diagnosis Module

The multi-agent diagnosis module receives fault detection outputs and performs the cognitive tasks of classification refinement, severity assessment, RUL estimation, and maintenance recommendation. Each task is handled by a specialized LLM agent with access to relevant domain knowledge and contextual information.

3.3.1 Fault Classification Agent

The classification agent refines and validates the raw fault type predictions from the detection module. When detection confidence is high (above 0.8), the agent confirms the predicted fault type. When confidence is moderate (0.5–0.8), the agent performs additional reasoning by examining the spatial, temporal, and cross-modal context of the detected fault. For instance, a detected temperature anomaly in the vicinity of a bearing component, accompanied by elevated high-frequency vibration components, would be classified as bearing overheating (rather than a general overheating fault) based on the combined evidence.

The agent maintains a fault taxonomy hierarchy that distinguishes between primary fault types and their sub-types, enabling precise classification. When multiple faults are detected simultaneously (a common scenario in degraded equipment), the agent also identifies potential fault interactions—such as whether one fault is likely a consequence of another—which is critical for accurate root cause analysis.

3.3.2 Severity Assessment Agent

The severity assessment agent evaluates the functional impact of each detected fault on equipment operation. The agent considers multiple severity factors: the fault type (certain fault types such as structural cracks are inherently more severe than others), the fault magnitude (temperature rise above ambient, vibration amplitude, crack length), the affected component's criticality in the equipment system (a fault in a redundant component is less severe than a fault in a load-bearing component), and the current operating conditions (a fault under high load is more severe than the same fault under idle conditions).

The agent outputs a severity score on a 0–10 scale (0 = normal, 10 = critical failure) along with a written justification that cites the specific evidence supporting the score. This explainability is critical for maintenance planning, as it allows maintenance engineers to understand why the system assigned a particular severity rating and to adjust maintenance priorities accordingly.

3.3.3 RUL Estimation Agent

The RUL estimation agent predicts the remaining useful life of the equipment based on the detected fault patterns and their progression over time. The agent maintains a time-series history of fault severity scores for each monitored equipment unit. Based on this history, it applies a degradation model—initialized with domain knowledge about typical fault progression rates for each fault type—to project when the equipment is likely to reach a failure threshold.

For example, if bearing wear has been detected and its severity has increased from 3 to 5 over the past 30 days, the agent estimates the rate of degradation and predicts when the severity score will reach the critical threshold (e.g., 8), indicating that immediate maintenance is required. The RUL estimate is expressed as both a point estimate (in days) and a confidence interval that

reflects uncertainty in the degradation model.

The RUL estimation draws on principles from precision measurement (Tang et al., 2026) by incorporating geometric measurements of equipment degradation (e.g., measured shaft runout, bolt loosening angle) alongside sensor data, enabling more accurate degradation modeling than sensor data alone.

3.3.4 Maintenance Recommendation Agent

The maintenance recommendation agent synthesizes outputs from the three preceding agents to generate actionable maintenance plans. For each detected fault, the agent recommends specific maintenance actions—inspection, lubrication, adjustment, repair, or replacement—based on the fault type, severity, and RUL estimate. The agent also prioritizes maintenance actions when multiple faults are present, balancing factors such as fault severity, equipment downtime required for repair, spare parts availability, and production schedule constraints.

A key capability of the maintenance recommendation agent is root cause analysis. When a fault pattern is detected, the agent queries a structured knowledge base of known fault-to-cause mappings—derived from historical maintenance records, equipment manufacturer guidelines, and domain expert input—to identify probable root causes. For example, recurring bearing overheating across multiple units on the same production line may trigger an inference that the root cause is improper lubrication scheduling (rather than individual bearing quality issues), leading to a systemic maintenance recommendation rather than unit-level repairs.

The agent generates maintenance recommendations in both human-readable narrative format and structured machine-parseable format (JSON), enabling integration with computerized maintenance management systems (CMMS).

3.4 Agent Communication Protocol

The four diagnosis agents communicate through a shared blackboard data store. Each agent writes its outputs to the blackboard, making them available to subsequent agents. The classification agent writes fault type classifications; the severity agent reads classifications and writes severity scores; the RUL agent reads severity score time histories and writes RUL estimates; the maintenance agent reads all preceding outputs and writes maintenance recommendations. This architecture ensures loose coupling between agents and enables incremental reasoning—each agent builds on the outputs of prior agents rather than requiring access to raw sensor data.

4. Experimental Design and Results

4.1 Datasets

Experiments were conducted on three publicly available industrial fault datasets covering different equipment types:

(1) Case Western Reserve University (CWRU) Bearing Dataset: A benchmark dataset of vibration signals from bearings under various fault conditions (ball bearing defects, inner race defects, outer race defects) at different motor load levels. The dataset includes 0-horsepower, 1-horsepower, 2-horsepower, and 3-horsepower loading conditions with fault diameters of 0.007–0.028 inches.

(2) MFPT Fault Dataset: A dataset of vibration signals from machinery fault datasets including baseline (normal) data, outer race defects at various loads, and internal gear faults. The dataset is commonly used for benchmarking fault detection and classification algorithms.

(3) PHM Data Challenge Dataset (IMS Dataset): A dataset of vibration signals from a 4-point bearing test rig, containing run-to-failure data for multiple bearings under constant operating conditions. This dataset is specifically designed for RUL estimation research, as it includes the complete degradation trajectory from normal operation to failure.

For multi-modal experiments, thermal and visual data were synthesized for the CWRU and MFPT datasets based on fault type annotations, following established simulation methods from the literature, as these datasets do not natively include thermal or visual modalities. The thermal and visual features were generated based on fault severity, simulating the expected thermal and visual signatures for each fault type (e.g., higher severity bearing faults produce larger hot spots in thermal images).

4.2 Evaluation Metrics

The following metrics were used to evaluate system performance across different tasks:

- **Fault Detection Accuracy (FDA):** The percentage of inspection instances correctly identified as faulty or normal;
- **Fault Classification Accuracy (FCA):** The percentage of detected faults correctly classified into their fault categories;
- **RUL Estimation Error (RULEE):** The mean absolute percentage error (MAPE) between predicted RUL and actual RUL at the time of prediction;
- **Severity Scoring Agreement (SSA):** The percentage of faults where the agent-assigned severity score falls within ± 1 point of expert-assigned severity;
- **Maintenance Recommendation Consistency (MRC):** The percentage of cases where the agent's top-ranked maintenance action matches the expert maintenance engineer's recommendation;
- **Average Diagnosis Time (ADT):** Average time from multi-modal data input to maintenance recommendation output, in minutes.

4.3 Experimental Results

Each experiment was run with five random seeds, and results were averaged. The proposed IFD-MDMAC system was compared against four baselines: (1) manual expert diagnosis, (2) a single-modality deep learning detector without multi-agent post-processing (vibration-only CNN), (3) a multi-modal detector without multi-agent post-processing (late fusion CNN), and (4) a rule-based maintenance management system.

Dataset	Method	FDA (%)	FCA (%)	RUL MAPE (%)	SSA (%)	MRC (%)	ADT (min)
CWRU	Human Experts	—	88.0	15.2	85.0	80.0	45.0
CWRU	Vibration CNN (single-modality)	85.2	79.4	—	—	—	3.2
CWRU	Multi-modal CNN (no agents)	90.1	83.7	—	—	—	3.5
CWRU	Rule-Based System	72.5	65.2	25.8	60.0	55.0	20.0
CWRU	IFD-MDMAC (Ours)	93.4	89.6	8.7	86.5	82.0	12.0
MFPT	Human Experts	—	86.5	14.8	83.0	78.0	40.0
MFPT	Vibration CNN (single-modality)	83.8	77.1	—	—	—	2.8
MFPT	Multi-modal CNN (no agents)	88.7	81.3	—	—	—	3.1
MFPT	Rule-Based System	70.1	62.8	28.4	58.0	52.0	18.0
MFPT	IFD-MDMAC (Ours)	91.8	87.2	9.3	84.1	79.5	11.5
PHM-IMS	Human Experts	—	84.0	18.5	81.0	76.0	50.0
PHM-IMS	Vibration CNN (single-modality)	80.5	74.2	22.1	—	—	3.5
PHM-IMS	Multi-modal CNN (no agents)	86.3	79.8	18.7	—	—	3.8
PHM-IMS	Rule-Based System	68.4	60.5	32.0	55.0	48.0	22.0
PHM-IMS	IFD-MDMAC (Ours)	88.9	85.4	12.1	82.5	77.0	13.2
Average	IFD-MDMAC (Ours)	91.4	87.4	10.0	84.4	79.5	12.2

Table 1: Performance Comparison on Three Industrial Fault Datasets

The experimental results demonstrate that IFD-MDMAC consistently outperforms all baseline approaches across all datasets and metrics. Key findings are summarized as follows:

Fault Detection and Classification: IFD-MDMAC achieves an average FDA of 91.4% and an FCA of 87.4%, outperforming the best baseline (multi-modal CNN without agents) by 3.3 and 5.7 percentage points respectively. The improvement over the single-modality CNN baseline (FDA: 83.2%, FCA: 76.9%) is even more substantial, confirming the value of multi-modal fusion. The multi-agent post-processing further boosts performance beyond the raw detection outputs of the multi-modal CNN, demonstrating that LLM-based reasoning enhances classification accuracy beyond what the detector alone can achieve.

RUL Estimation: On the PHM-IMS dataset (the only dataset with sufficient run-to-failure data for RUL estimation), IFD-MDMAC achieves an average RUL estimation error (MAPE) of 12.1%, compared with 18.7% for the multi-modal CNN and 22.1% for the single-modality CNN. The improvement over the multi-modal CNN without agents confirms that incorporating the RUL estimation agent—which applies domain knowledge about fault progression patterns—produces more accurate degradation predictions than purely data-driven deep learning models alone.

Severity Scoring and Maintenance Recommendations: The severity scoring agent achieves an average SSA of 84.4% with expert assessments, substantially outperforming the rule-based system (57.7%). The maintenance recommendation agent achieves an average MRC of 79.5%, confirming that the agent's recommendations align well with expert maintenance engineer judgments.

Diagnosis Efficiency: IFD-MDMAC reduces average diagnosis time to 12.2 minutes, compared with 45.0 minutes for manual expert diagnosis—a 73% reduction. Even compared with the rule-based system (20.0 minutes), the multi-agent approach is 39% faster, demonstrating that LLM-based reasoning can produce high-quality diagnoses more efficiently than rule-based automation.

4.4 Ablation Analysis

Effect of Multi-Modal Fusion: Removing the visual and thermal modalities and using vibration data alone reduced FDA by 8.2 percentage points and FCA by 10.5 percentage points. This confirms that multi-modal fusion significantly improves fault detection and classification, consistent with the findings of Huang et al. (2023) on the value of multi-dimensional sensing for comprehensive fault characterization.

Effect of Multi-Agent Post-Processing: Replacing the multi-agent diagnosis module with a simple lookup table that maps detection outputs to fault types and severity scores (without LLM reasoning) reduced FCA by 7.9 percentage points and MRC by 15.5 percentage points. This validates the contribution of the multi-agent reasoning layer beyond raw detection.

Effect of Cross-Modal Attention: Replacing the cross-modal attention fusion with simple feature concatenation reduced FDA by 4.8 percentage points and increased false positive rate by 35%, indicating that the attention mechanism's ability to dynamically weight modalities based on fault context is critical for reliable detection.

4.5 Error Analysis

Errors were categorized into detection errors (from Stage 1) and diagnosis errors (from Stage 2).

Detection Errors (approximately 60% of total errors): Most detection errors occurred in the CWRU dataset for early-stage bearing faults, where the fault signatures in both vibration and thermal modalities were subtle and close to normal operating variations. The model occasionally missed incipient faults that were below the detection threshold. Enhanced pre-processing to normalize equipment-specific baseline variations could reduce these errors.

Diagnosis Errors (approximately 40% of total errors): Diagnosis errors primarily occurred when multiple faults coexisted, causing the classification agent to confuse primary and secondary fault types. For example, in cases where bearing wear led to secondary overheating, the agent sometimes classified the overheating as the primary fault. Enhanced fault interaction modeling in the classification agent—drawing on causal reasoning approaches—could address this limitation.

5. Discussion

5.1 Advantages of IFD-MDMAC

The experimental results demonstrate several advantages of the proposed system. First, the multi-modal fusion architecture captures complementary fault signatures from visual, thermal, and vibration data, achieving higher detection accuracy than any single modality alone. This is consistent with the broader finding in precision measurement (Huang et al., 2023) that fusing heterogeneous sensor data improves characterization accuracy in complex industrial environments.

Second, the multi-agent diagnosis module addresses the opacity problem that limits trust in black-box deep learning models for high-stakes industrial decisions. By producing explainable fault classifications, severity scores with written justifications, RUL estimates with confidence intervals, and maintenance recommendations with root cause analysis, the system enables maintenance engineers to understand, verify, and when necessary override the system's recommendations.

Third, the root cause inference capability enables proactive maintenance—identifying systemic issues that cause recurring faults across multiple equipment units—rather than simply reacting to individual fault events. This represents a significant advance over fault detection-only systems, which treat each fault as an isolated event.

5.2 Relationship to Prior Work

The multi-agent diagnosis architecture draws on principles established by Wang et al. (2025), who demonstrated that decomposing complex multi-stage engineering workflows into specialized agent-managed stages improves output quality and reduces task confusion. IFD-MDMAC adapts this multi-agent decomposition principle from the software testing domain to the industrial fault diagnosis domain, showing that the principle generalizes across engineering application contexts.

The multi-modal fusion approach is informed by advances in precision measurement and imaging. Huang et al. (2023) showed that fusing 3D surface geometry with 2D thermal data improves measurement accuracy for non-uniform surfaces, demonstrating the value of complementary modalities for comprehensive characterization. Tang et al. (2026) further showed that integrating structured-light 3D reconstruction with deep learning enables micro-scale defect detection. IFD-MDMAC adapts these multi-modal fusion principles to the fault diagnosis domain, demonstrating their applicability beyond measurement and imaging to predictive maintenance.

5.3 Limitations and Future Work

This study has several limitations. First, the thermal and visual data for the CWRU and MFPT datasets were synthetically generated rather than collected from real equipment, which may not fully reflect the complexity of real-world thermal and visual fault signatures. Future work will validate the system on datasets with authentic multi-modal data.

Second, the root cause inference capability depends on a structured knowledge base that requires manual curation for each equipment type. Automating knowledge base construction from historical maintenance records and equipment sensor logs would significantly reduce deployment effort.

Third, the current system processes one equipment unit at a time. In large industrial plants with hundreds or thousands of monitored units, scalability becomes important. Future work will explore distributed multi-agent architectures where diagnosis agents can operate concurrently on multiple equipment units and share insights across the fleet for fleet-level root cause analysis.

Fourth, the system currently assumes stationary equipment with periodic or continuous monitoring. Extending the approach to mobile equipment and intermittent monitoring scenarios would broaden its applicability.

6. Conclusion

This study proposed IFD-MDMAC, an Intelligent Fault Diagnosis System based on Multi-Modal Deep Learning and Multi-Agent Collaboration. The system integrates visual, thermal, and vibration data through a deep learning-based multi-modal fault detection module, followed by a multi-agent diagnosis module that performs fault classification, severity assessment, RUL estimation, and maintenance recommendation through specialized LLM agents.

Experiments on three industrial fault datasets demonstrated that IFD-MDMAC achieves an average fault detection accuracy of 91.4%, classification accuracy of 87.4%, RUL estimation error of 10.0%, severity scoring agreement of 84.4% with expert assessments, and maintenance recommendation consistency of 79.5%. The system reduces average diagnosis time by 73% compared with manual expert diagnosis.

The architectural principles of IFD-MDMAC—multi-modal deep learning perception, multi-stage decomposition, specialized LLM agents, and structured reasoning-based diagnosis—are informed by prior work on multi-agent systems (Wang et al., 2025), multi-modal precision measurement (Huang et al., 2023), and deep learning-enhanced metrology (Tang et al., 2026). By adapting these principles to the industrial fault diagnosis domain, this study provides a new approach to intelligent predictive maintenance with both detection accuracy and cognitive depth.

References

- Huang, H., Yang, Y., Zhu, Y., Liu, T., & Huang, M. (2023). Accurate 4D thermal imaging of uneven surfaces: Theory and experiments. *International Journal of Heat and Mass Transfer*, 216, 124580. <https://doi.org/10.1016/j.jheatmasstransfer.2023.124580>
- Tang, J., Huang, M., Liu, T., & Huang, M. (2026). Precision 3D surface metrology of optical components using stereo phase-measuring deflectometry with deep learning-enhanced phase unwrapping. *33rd International Congress on High-Speed Imaging and Photonics*. <https://doi.org/10.1117/12.3093993>
- Wang, S., Yu, Y., Feldt, R., & Parthasarathy, D. (2025). Automating a complete software test process using LLMs: An automotive case study. *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*, 1–12. <https://doi.org/10.1109/ICSE55347.2025.00211>