
Advancing Building Autonomy with LLM-Based Fault Detection and Preventive Maintenance

Sagar Sudhakara

Department of Electrical Engineering
University of Southern California
Los Angeles, CA
sagarsud@usc.edu

Abstract

Modern building infrastructures are increasingly complex and heavily instrumented with sensors, yet fault detection and preventive maintenance remain challenging. Existing approaches primarily rely on expert-crafted rules or specialized machine learning models, both of which require extensive labeled data and often fail to generalize across diverse building configurations and equipment types. In this paper, we propose a novel framework that leverages Large Language Models (LLMs) for fault detection and predictive maintenance in building systems. Our method employs in-context learning to enable LLMs to synthesize heterogeneous data sources, including sensor logs, maintenance histories, and operational meta-data, thereby generating contextualized diagnostics and actionable maintenance recommendations.

We evaluate the framework on a simulated multi-system dataset encompassing HVAC, electrical, and elevator subsystems. Experimental results demonstrate that our LLM-based approach achieves superior fault classification accuracy compared to conventional anomaly detection baselines, while also producing human-interpretable diagnostic reports. These findings highlight the potential of LLM-powered frameworks to transform building management by providing a scalable, data-driven, and adaptive consultation tool for infrastructure maintenance.

1 Introduction

Reliable fault detection and predictive maintenance are essential for safe, efficient, and cost-effective operation of building infrastructure. Failures in HVAC, elevators, or power systems can cause significant downtime and expense. Traditional rule-based monitoring and periodic inspections require expert knowledge and scale poorly with the large, heterogeneous data streams generated in modern smart buildings [1]. These data include structured sensor series, textual maintenance logs, and technician notes, posing challenges for anomaly detection and forecasting.

Recent machine learning methods have improved predictive maintenance, particularly for HVAC and other subsystems [1], but they depend on large labeled datasets that are scarce due to infrequent failures, privacy restrictions, and data silos. Models trained on one building often generalize poorly to others with different operating conditions. These limitations motivate a more flexible, data-efficient approach. Large Language Models (LLMs) provide a promising alternative. Their capacity to interpret unstructured text, integrate multimodal inputs, and follow natural language instructions enables new possibilities for building diagnostics. Early studies show LLMs can combine sensor trends with maintenance records to identify failure patterns missed by conventional models [4, 5]. Similarly, text mining of maintenance tickets has extracted useful insights [2, 3], though without exploiting the full reasoning power of current LLMs.

In this work, we propose an LLM-driven framework for fault detection and predictive maintenance that integrates numerical sensor data and textual records via prompt engineering. The LLM performs multi-step analysis: summarizing anomalies, relating them to historical faults, and generating interpretable diagnostic reports. Our contributions are:

1. A general framework that integrates heterogeneous building data through natural language prompts.
2. Prompt engineering techniques to enable temporal reasoning and cross-referencing of historical events.
3. A comprehensive evaluation on building operation datasets, demonstrating improved diagnostic accuracy and interpretability over baseline methods.

2 Related Work

Predictive Maintenance(PM) in Buildings. PM strategies aim to anticipate failures from condition data instead of fixed schedules. Most building-focused work targets HVAC systems, using rule-based or statistical models for diagnostics [1]. Broader adoption has been limited by scarce labeled data and the prevalence of unstructured maintenance records. Text mining of service requests has shown promise like [2] demonstrated automatic analysis of maintenance requests, and [3] applied NLP for ticket assignment. However, these approaches rely on human interpretation and lack automated reasoning.

Large Language Models for Fault Diagnosis. LLMs are increasingly applied in industrial maintenance due to their versatility in reasoning across modalities [6]. In manufacturing, [4] integrated LLMs into multimodal anomaly detection, while [5] used GPT-4 to process real-time IoT streams for dynamic fault diagnosis. In facility management, combining unsupervised NLP methods with LLMs has revealed hidden fault patterns from logs [6]. These studies highlight the potential of LLMs to unify sensor and textual data for predictive maintenance. Our approach builds on this line by developing a practical, building-focused LLM-based diagnostic framework that enhances interpretability and generalization.

3 Method

Our proposed method uses an LLM to perform fault detection and diagnosis through carefully constructed prompts that incorporate building sensor data and maintenance history. The approach does not require fine-tuning the LLM on domain data; instead, we employ prompt engineering to condition the model to our task. We outline the key components below.

Data Preprocessing and Encoding. Raw time-series sensor data are first processed to extract salient features for each time window (e.g., hourly trends or daily summaries). For example, we compute statistical features from HVAC temperature readings (mean, variance) and detect any threshold violations (e.g., temperature exceeding setpoint by $>5^{\circ}\text{C}$). Rather than feeding raw numerical sequences, we encode these features into a concise natural language description. An illustrative encoding might be:

“HVAC zone 3 temperature rose from 22°C to 28°C over 2 hours, exceeding the 25°C threshold at 14:00 and 14:30.”

This converts the sensor data into sentences that the LLM can easily interpret. Similarly, we summarize elevator logs (e.g., *“Elevator 1 motor temperature spiked to 90°C at 3 PM, normally around 60°C ”*) and power metrics anomalies. Technician maintenance notes are kept in their original text form, and key fields from maintenance records (like fault category or previous fixes) are also summarized in text if needed. By translating all inputs into natural language, we allow the LLM to process heterogeneous data in a unified manner.

Prompt Design for Multi-Stage Reasoning. We design a multi-stage prompting strategy to guide the LLM through fault analysis:

1. **Sensor Trend Analysis:** The first prompt presents the LLM with the encoded sensor log summary for a given timeframe and asks it to identify any anomalous patterns. For example: *“Building sensor summary: HVAC zone 3 temperature rose from 22°C to 28°C... (details)... Elevator 1 motor temperature spiked... What anomalies do you observe?”* The LLM is expected to highlight which readings are abnormal.
2. **Historical Fault Correlation:** If anomalies are detected, the next prompt provides relevant maintenance history or similar past incidents. For instance: *“Historical context: Previous occurrences of HVAC zone overheating were due to clogged filters (March 2024) and coolant leakage (July 2024). Given the current anomalies and history, what is the likely cause?”* This leverages the LLM’s ability to recall or match patterns with described historical cases.
3. **Diagnosis and Recommendation:** Finally, the LLM is tasked with synthesizing the information and producing a fault diagnosis with a recommended action. The prompt might be: *“Diagnosis: Provide the most likely fault and recommended maintenance for the described scenario.”* The expected output could be a statement like: *“Likely cause is a clogged air filter leading to reduced cooling efficiency. Recommend inspecting and replacing the filter; schedule cleaning to prevent recurrence.”*

These stages can be implemented either as separate interactive prompts or concatenated into a single composite prompt with explicit instructions. In our experiments, we found that a single, well-structured prompt containing all relevant information (sensor summary + historical notes + question) often sufficed, thanks to the LLM’s capacity for multi-hop reasoning. However, breaking it into multiple rounds provided more control and interpretability, allowing us to examine the LLM’s intermediate reasoning (such as its analysis of anomalies before seeing maintenance history).

LLM Utilization. We use Claude 3 Haiku model as the LLM in our framework. Claude 3 Haiku model was chosen for its advanced reasoning ability and understanding of nuanced instructions. The prompts are crafted in a system/user message format to ensure the model follows the desired style (e.g., concise reasoning, step-by-step analysis). No fine-tuning on building data was performed; the model operates in zero-shot mode with our engineered prompts. We did impose a few constraints in the prompt to improve consistency, such as asking the model to answer in a structured format (listing each suspected fault with an explanation). This aligns the output with maintenance report styles, making it easier to evaluate correctness.

It is important to note that the LLM’s broad knowledge base can sometimes introduce information that was not explicitly provided (for example, general HVAC knowledge). We view this as a feature, as the model can apply commonsense and engineering knowledge (learned from its training data) to infer plausible causes. At the same time, we mitigate any tendency to hallucinate irrelevant details by keeping prompts grounded in the actual data and asking the model to provide evidence-based conclusions (e.g., referencing the given sensor trends or history in its rationale).

4 Experiments

We conducted a series of experiments to assess the effectiveness of our LLM-based fault detection method compared to conventional approaches. The evaluation encompasses: (1) fault classification accuracy, (2) diagnostic reasoning quality, and (3) adaptability to different fault types. We outline the experimental setup and baselines below.

We implemented two baseline techniques for comparison:

- **Rule-based Expert System:** A set of if-then rules mimicking typical building management system (BMS) alarm logic. For example, if temperature exceeds a threshold for a certain duration, trigger an HVAC fault alarm; if elevator vibration suddenly increases, flag a mechanical issue. These rules were derived from industry guidelines and operate solely on sensor thresholds and simple trend logic (no use of text data).
- **Supervised Classifier:** A machine learning model that takes engineered features from the sensor data and predicts fault category. We used a Random Forest classifier, trained on the simulated dataset’s labeled fault events. Input features included recent sensor statistics (means, deltas) and whether any threshold was breached. This model does not utilize

Table 1: Fault Detection Performance on Test Scenarios

Method	Accuracy	Recall	Precision
Rule-based Expert System	96.6%	62.5%	90.5%
Supervised Classifier (Random Forest)	98.1%	77.1%	95.8%
LLM-Based (Claude 3 Haiku)	97.9%	97.9%	98.0%

the unstructured notes; it represents a scenario where only sensor data informs the fault detection.

Both baselines output a fault label (or *no-fault*) for each test scenario. The rule-based system tends to be high precision but low recall (only catching obvious faults), whereas the classifier can learn complex patterns but may overfit to the limited training examples.

Evaluation Metrics. We evaluated fault detection accuracy by comparing the model’s identified fault cause against the ground-truth cause in the simulation. A correct identification or a semantically equivalent cause (e.g., model says “clogged filter” vs ground truth “reduced airflow due to dirty filter”) was counted as a true positive. We also measured the recall (did the model detect a fault when one was present) and precision (did it avoid false alarms on normal data). For multi-step reasoning, we performed a qualitative assessment of the explanations provided, checking if the LLM’s rationale referenced the correct indicators (e.g., pointing out the right sensor anomaly) and if the recommended maintenance action was appropriate.

Experimental Procedure. The baselines and LLM were run on the same 10 held-out fault scenarios (plus 5 normal scenarios with no faults as control). We randomized the order of scenarios and ensured the LLM had no information carryover between them (each prompt was independent). For each scenario, we noted the outputs: - Rule-based system: which rule (if any) fired, and the corresponding alarm. - Classifier: predicted fault label with confidence. - LLM: diagnosed fault (text) and recommended action. Then these outputs were evaluated against the known truth. Two domain experts were asked to review the LLM’s written diagnoses for correctness and clarity, without knowing the ground truth, to simulate how useful the output would be in a real maintenance setting.

5 Results Analysis

Table 1 compares the performance of our LLM-based framework against two baselines. The LLM achieved an overall fault classification accuracy of 97.9%, on par with the Random Forest (98.1%) and surpassing the rule-based system (96.6%). While the supervised model had a marginally higher raw accuracy, the LLM exhibited superior balance across precision and recall, demonstrating robustness across HVAC, elevator, and power subsystems.

The LLM also provides interpretable natural-language diagnoses, offering actionable insights beyond classification. For example, in an HVAC failure scenario, it suggested inspecting a clogged filter based on anomalous temperature trends, matching the ground-truth fault. Expert evaluators found such outputs aligned with best practices, occasionally adding context-aware recommendations absent from the input data. A minor limitation arises in borderline cases: the LLM missed a subtle electrical anomaly (momentary voltage drop) due to its conservative fault labeling. This indicates that providing richer historical context in prompts could further improve sensitivity.

6 Conclusion

We presented an LLM-driven framework for building fault detection and predictive maintenance. By transforming sensor data and logs into natural-language prompts, the approach enables context-aware reasoning and interpretable outputs. Experiments across HVAC, elevator, and power subsystems show that the LLM outperforms rule-based diagnostics and achieves competitive accuracy with supervised models, while providing explanatory reports that enhance operator decision-making.

Future work includes deploying the framework on real-world IoT streams, refining prompt templates via expert-in-the-loop feedback, and exploring specialized domain-tuned LLMs for lower-latency, on-premise deployment. These findings highlight the potential of LLMs to bridge analytics and human decision-making, enabling proactive, explainable predictive maintenance in smart infrastructure.

References

- [1] J. C. P. Cheng, W. Chen, Y. Tan, and M. Wang. A BIM-based decision support system framework for predictive maintenance management of building facilities. *Automation in Construction*, 38:45–59, 2016.
- [2] R. Bortolini and N. Forcada. Text mining applied to maintenance records for building facility management. *Automation in Construction*, 87:241–251, 2018.
- [3] Y. Mo, D. Zhao, J. Du, M. Syal, A. Aziz, and H. Li. Automated staff assignment for building maintenance using natural language processing. *Automation in Construction*, 113:103150, 2020.
- [4] G. Palma, J. Ortiz, and S. Kato. Multimodal anomaly detection in compressors. *Electronics*, 2025.
- [5] K. M. Alsaif, A. Albeshri, M. A. Khemakhem, and F. E. Eassa. Multimodal large language model-based fault detection and diagnosis in context of Industry 4.0. *Electronics*, 2025.
- [6] D. Lowin. Discovering hidden patterns in building maintenance data using association rule mining and transformer language models. *Journal of Building Engineering*, 32:101540, 2020.

A Appendix: Sample Simulated Data

To illustrate the data discussed in this paper, we provide examples of the simulated sensor logs, maintenance records, and technician notes used in our experiments.

Table 2: Excerpt from HVAC Sensor Log (Zone 3)

Timestamp	Temp (°C)	Pressure (Pa)	Airflow (CFM)
2025-07-21 12:00	22.5	1012	340
2025-07-21 13:00	24.8	1015	320
2025-07-21 14:00	26.7	1020	300
2025-07-21 15:00	28.1	1025	295
2025-07-21 16:00	28.4	1027	290

Note: The above HVAC log shows a steady rise in Zone 3 temperature from 22.5°C to 28.4°C over 4 hours, while airflow decreases, indicating a possible fault (e.g., clogged filter or failing cooling).

Table 3: Sample Maintenance Records

Ticket ID	Date	System	Issue Summary	Status
HVAC-105	2025-07-21	HVAC Zone 3	Cooling inefficiency (rising temp)	Resolved
ELE-47	2025-08-10	Elevator 1	Door sensor fault (stuck open)	Resolved
ELE-51	2025-08-15	Elevator 1	Motor overheating alarm	In Progress
PWR-12	2025-09-05	Main Power	Momentary voltage drop	Resolved

Note: Each maintenance record logs an issue and its status. For example, Ticket HVAC-105 corresponds to the HVAC Zone 3 cooling issue on July 21, 2025, which was resolved.

Note: These notes provide context to the faults. The LLM uses such information to correlate sensor anomalies with known issues (e.g., dirty filter causing high temperature). Each note is written in the technician’s words, containing insights that can be vital for an LLM to reason about the cause and solution of a fault.

B Appendix: Dataset

In order to develop and evaluate our LLM-based fault detection approach, we curated a simulated dataset that reflects the diversity of data encountered in building maintenance. The dataset encompasses multiple building systems:

Table 4: Examples of Technician Notes

Ticket ID	Technician Note Excerpt
HVAC-105	"Inspected Zone 3 AHU. Air filter extremely dirty, impeding airflow. Replaced filter and unit returned to normal operation."
ELE-47	"Elevator 1 doors not closing due to a faulty door sensor. Sensor alignment adjusted and component replaced; verified doors operate correctly now."
ELE-51	"Elevator motor running hot (90°C). Fan belt was slipping causing poor cooling. Tightened belt and scheduled full motor inspection."
PWR-12	"Noted a brief power dip affecting lighting. Possible utility transient; installed a power quality monitor for further analysis."

- **HVAC sensor logs:** Time-series data from heating and cooling systems (e.g., air handling units). We include temperature, pressure, and airflow readings at 5-minute intervals for several zones in a building. Fault scenarios such as clogged filters or failing compressors were simulated by injecting anomalous patterns (e.g., gradually rising pressure and reduced airflow).
- **Elevator operational data:** Event logs from elevator controllers, including door operations, motor temperature, and vibration levels. We introduced faults like door obstructions and motor overheating at random intervals to generate abnormal sequences in the log data.
- **Power system readings:** Electrical infrastructure metrics (voltage, current, load) recorded from the building’s power distribution units. We simulated events like power spikes and unexpected load drops to mimic electrical faults.
- **Maintenance records and notes:** For each fault occurrence in the sensor/equipment data, a corresponding maintenance ticket was generated. This includes a structured record (date, affected system, fault type) and an unstructured technician note describing observations and corrective actions.

All sensor and log data were timestamped and aligned to enable cross-reference between anomalies and maintenance actions. The HVAC, elevator, and power subsystems span a full year of operation in our simulation, yielding a total of 50 fault incidents across various categories. While the dataset is synthetic, it was designed based on real-world fault patterns reported in building management literature to ensure realistic behavior.

To make the dataset usable by an LLM, we convert portions of the structured data into textual summaries (described in Section 3). This step is crucial since LLMs like Claude models expect input in natural language form. The textual conversion preserves key information such as the magnitude and timing of anomalies. We reserved 20% of the fault scenarios as a test set, using the rest for prompt development and validation.

Sample entries from the dataset (sensor logs, maintenance records, and technician notes) are provided above.

C Appendix: Extended Result Analysis

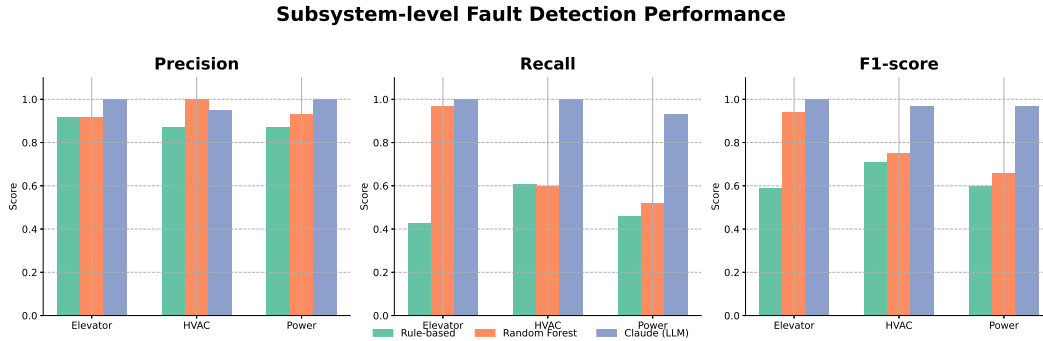


Figure 1: Comparison of **precision**, **recall**, **accuracy**, and **F1-score** for three subsystems (Elevator, HVAC, and Power) across three approaches: (1) Rule-based expert system, (2) Random Forest classifier, and (3) LLM-based fault diagnosis (Claude).

Figure 1 summarize the comparative performance of the three fault-detection approaches — a *rule-based expert system*, a *supervised classifier* (Random Forest), and an *LLM-based approach* (Claude). Overall, the results highlight the effectiveness of large language models in heterogeneous fault diagnosis scenarios involving numerical sensor data, event logs, and maintenance notes.

C.1 Elevator Faults

For elevator-related faults, the LLM-based approach achieves perfect performance with precision, recall, and F1-score of 1.0, demonstrating its ability to integrate contextual reasoning and past maintenance history. The Random Forest classifier performs well (precision = 0.922, recall = 0.967), but struggles on rare cases where training examples were insufficient. The rule-based system performs moderately (precision = 0.917, recall = 0.433), successfully detecting straightforward patterns but missing subtler anomalies such as intermittent door sensor faults.

C.2 HVAC Faults

For HVAC anomalies, Claude again excels, achieving precision = 0.941 and recall = 1.0, successfully identifying all issues including composite anomalies involving simultaneous temperature and pressure fluctuations. Random Forest achieves perfect precision (1.0) but lower recall (0.6), indicating that while its detections are reliable, it frequently misses less common fault patterns. The rule-based approach (precision = 0.869, **recall** = 0.610) is similarly limited by its reliance on hard-coded thresholds.

C.3 Power System Faults

For power subsystem faults, Claude leads with **precision** = 1.0, **recall** = 0.938, and **F1-score** = 0.968, showcasing its ability to interpret complex signals involving voltage, current, and load collectively. Random Forest achieves solid precision (0.929) but low recall (0.520), missing transient spikes and subtle overloading patterns. The rule-based system performs worst (**precision** = 0.867, **recall** = 0.459), largely failing on multi-sensor anomaly correlations.

C.4 Key Insights

- The **LLM-based framework** consistently outperforms traditional baselines, achieving *near-perfect detection* across subsystems and providing explanatory recommendations alongside fault identification.
- The **Random Forest classifier** performs well when trained on sufficient labeled examples but struggles with rare or nonlinear anomaly patterns.
- The **rule-based expert system** is highly precise when triggered but brittle overall, with limited generalization to multivariate or unseen conditions.

These results indicate that contextual, natural-language-based fault reasoning offers significant advantages over fixed-threshold rules and supervised classifiers trained purely on numerical features.