# AUGMENTING INDUSTRIAL MAINTENANCE WITH LLMS: A BENCHMARK, ANALYSIS, AND GENERALIZATION STUDY

#### **Anonymous authors**

000

001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027 028 029

031

032

033

034

037

040

041

042

043

044

045

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Monitoring the life cycle of complex industrial systems often relies on expertly curated temporal conditions derived from sensor data, a process that requires significant time investment and deep domain expertise. We explore the potential of utilizing Large Language Models (LLMs) to generate context-aware and accurate recommendations for maintenance based on their ability to reason and generalize on temporal sensor conditions. To this end, we formulate a novel pipeline that systematically converts human-authored symbolic conditions into a multiple-choice question answer (MCQA) dataset. We apply our pipeline by creating DiagnosticIQ, a 6,000+ MCQA dataset covering 16 different types of physical assets that represent real-world maintenance use cases. We assess 15 state-of-the-art Large Language Models (LLMs) with this dataset and create a leaderboard for the maintenance action recommendation task. Furthermore, we evaluate and demonstrate the practical utility of DiagnosticIQ in two key aspects. First, as a knowledge base to enhance maintenance action recommendations, and secondly, as a fine-tuning resource to fine-tune a specialized LLM that generalizes across previously unseen assets to facilitate the rule creation process.

#### 1 Introduction

Industrial complex equipments such as wind turbines, air handling units, and chillers require significant domain expertise for appropriate and effective operations, maintenance and tuning. These equipments are frequently deployed in operationally critical environments, such as health care organizations and large data centers where enhancing operational reliability and efficiency are critical. To achieve this, many Industrial facilities have integrated automated monitoring systems such as Internet of Things (IoT) solutions which continuously capture sensor data reflecting the operational state of the equipment and the interconnected elements of the equipments. While these systems can detect anomalies by monitoring predefined conditions, they generally provide limited guidance on appropriate corrective actions once issues are identified.

Consider Bob, a facility manager responsible for maintaining HVAC systems in a data center. Bob configures rule-based alarms by analyzing sensor data and asset metadata, defining asset-specific logical conditions such as Temperature  $> 80^{\circ}$ F or Enthalpy < 15 BTU/lb

Condition 1 Condition 2 that indicate potential faults (Figure 1). When these conditions are met, alerts are automatically triggered and sent to Alice, an operator, enabling timely and coordinated actions, such as inspecting broken valves, slipping, or misaligned belts.

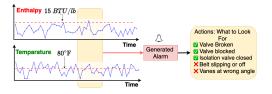


Figure 1: Faults to Fixes: Operation Workflow from Monitoring IoT Stream to Alarm Generation to Actionable Maintenance Recommendations

At the center of this workflow lies the configuration of recommended actions or inspection steps deciding what to look for or what maintenance should be performed once an alert is triggered. Despite timely notifications of abnormal conditions, determining the specific maintenance, repair,

or verification steps often exceeds Bob's expertise. This challenge is magnified in operational environments where there is a variety of assets from different manufacturers, with different operating modes, equipped with hundreds of sensors, resulting in a vast number of conditions to monitor simultaneously. For each abnormal condition, identifying what to inspect or repair requires specialized knowledge of asset-specific failure modes and mechanical systems expertise typically gained through years of hands-on experience. Can LLMs help bridge this gap?

Recognizing this challenge, intelligent recommendation systems that translate complex sensor data into actionable maintenance steps are critical for effective industrial asset management. Use of LLMs in discovering rules in an automated manner from labeled operational data is demonstrated recently Zhang et al. (2025), and motivated in recent survey articles Raza et al. (2025); Su et al. (2024). Clearly, this stream of work will help in industrial applications such as predictive maintenance and signal monitoring Cook et al. (2019); Kanawaday & Sane (2017); Beghi et al. (2016); Shah & Tiwari (2018). However, the key step of connecting these discovered rules to actionable guidance for technicians remains unaddressed. Recent advances in language models offer promise for this action recommendations task Zhong et al. (2024), determining the correct maintenance or repair actions following alarms but their systematic evaluation is hindered by a lack of realistic, standardized benchmarks.

To address this gap, we present DiagnosticIQ, a novel benchmark suite for industrial maintenance action recommendation. Grounded in real-world scenarios, this suite features a primary multi-choice question-answering dataset along with several variants, each designed to systematically evaluate a specific capability vital for LLMs in this domain, such as reasoning, generalization, and robustness. We further analyze a set of frontier models under these axes to establish a strong baseline revealing the current strengths and limitations of LLMs under this domain.

#### Our contributions are as follows:

- We formalize and implement a novel deterministic dataset generation pipeline that converts expert-authored symbolic rules into MCQA dataset
- We release DiagnosticIQ and its specialized variants, a first-of-its-kind benchmark dataset with about 6,690 MCQA, expertly validated, based on 120 operational rule-action pairs.
- We benchmark over 15 large language models (including Claude, Gemini, GPT variants) and establish a Maintenance Action Recommendation Leaderboard to foster community evaluation and progress.
- We systematically evaluate LLMs underaxes of Reasoning, Generalization, Robustness and demonstrate the utility of DiagnosticIQ as a external knowledge base as well as a finetuning resource for the task of maintenance action recommendation.

#### 2 Related Work

Building QA datasets in specialized domains has been an emerging trend across the board such as telecommunications Lee et al. (2024), climate Schimanski et al. (2024), finance Chen et al. (2024), healthcare Ray et al. (2024); Sviridova et al. (2024), IT operations Zhang et al. (2024a), power plants Hong et al. (2024), and scientific disciplines Bhattacharjee et al. (2024). In this section, we review the most relevant papers related to our work, with a particular focus on multi-choice QA, statistical and fine-tuning methods, rule generation and the role of domain experts.

Multi-Choice Question Answering (MCQA) has become a popular construct in the LLM world due to its ease of evaluation and closed form. This is evident from dozens of recent works such as **TruthfulQA** Lin et al. (2022), **GPQA** Rein et al. (2023), **MMLUPro** Wang et al. (2024), **Failure-SensorIQ** Constantinides et al. (2025), **CureBench** Gao (2025), and **Multi-Modal AD** Jiang et al. (2025). These datasets are designed to evaluate various aspects of LLMs, such as common-sense reasoning, domain understanding, and multimodal context. The key difference lies in how these questions are generated in the first place and then validated by a domain expert if needed. We observed that the community has paid less attention to cases where part of the question is presented as a rule. **ComplexBench** Wen et al. (2024) is closest to our work in term of instruction, as it highlights the importance of evaluating LLMs on their ability to follow complex instructions such as And, Selection, and Chain. However, it does not explicitly cover rules. As discussed in Zhang et al. (2025),

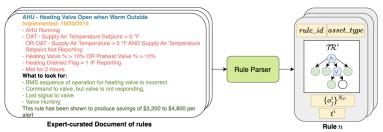


Figure 2: Example expert-curated rules and associated conditions/metadata for dataset construction.

rules encode domain knowledge effectively and benefit LLMs, and interestingly, many industrial monitoring systems utilize human-crafted rules for monitoring tasks. Therefore, we focus on building an MCQA dataset where the question part contains a rule.

Given an answer as part of the choices in MCQA, an LLM model may make a random selection and still achieve a favorable result. This limitation has motivated the research community to develop innovative statistical solutions and tools, such as validation tests Zheng et al. (2024); Zhang et al. (2024c); Robinson et al. (2023), PertEval Li et al. (2024a); Ye et al. (2024), and other approaches Zhu et al. (2024). The development of new, unseen MCQA datasets has further benefited the community. Furthermore, MMLUPro has adopted a method of increasing the number of options from four to ten and demonstrated that this adjustment has a non-trivial impact. However, constructing a robust solution for such complex questions remains an open challenge, necessitating innovations such as the adoption of a recommendation-based module, as discussed in this paper.

Statistical analysis of MCQA datasets still requires an additional level of human or domain expert validation, particularly for mission-critical applications. In most domain-specific QA datasets, experts are employed to participate in the quiz to quantify the difficulty of the prepared questions. Very limited analysis has been conducted on truly exercising the validation of the LLM's reasoning abilities, such as in the medical domain Sviridova et al. (2024). Evaluating the generated rationales/reasoning not only ensures correct answers but also builds confidence in the model's outputs. Unlike most prior studies, we examine whether LLMs generate accurate rationales alongside correct answers, evaluating explanations against end-user needs rather than relying solely on accuracy.

Operational rules are common in industrial domains, as demonstrated by Oracle's Maintenance Cloud Service Oracle (2025), but they require significant expert involvement Zhang et al. (2025). For many enterprise customers, smaller language models will be key, as they provide a practical way to embed domain-specific knowledge directly into the model. Methods such as Supervised Fine-Tuning (SFT) and Group Relative Policy Optimization (GRPO) Shao et al. (2024) have become key methods of this process, making generalization tests essential. Industrial settings often demand transferring rules written for one physical asset to another, further complicating the task. Despite the importance of these challenges, systematic comparisons in Industry 4.0 contexts remain limited, motivating our cross-asset knowledge transfer evaluation.

#### 3 SYMBOLIC CONDITIONS TO MCQA

This section details our methodology, beginning with the introduction of our pipeline that systematically converts symbolic, human-authored rules into a MCQA format. We then describe its application in the domain of industrial asset maintenance to create our benchmark dataset, DiagnosticIQ. The generation pipeline consists of three primary stages: (1) parsing rule documents into structured representations, (2) converting these representations into Disjunctive Normal Form (DNF), and (3) Selecting sets of actions. Finally, we discuss the development of several variants of DiagnosticIQ, each tailored to evaluate a specific capability of LLMs such as their temporal reasoning and generalization that is critical for the maintenance recommendation task.

#### 3.1 Input: Expert-Curated Rule Documents

The rule documents originated from the Smarter Buildings initiative, where Reliability Engineers, System Administrators, and a Rules Logic committee collaboratively developed and iteratively refined fault-detection logic across multiple equipment types (e.g., Air Handlers, Chillers, Boilers) over several years to expand coverage and maintain diagnostic accuracy. A detailed description of the

rule development process is provided in Appendix B. The pipeline generation process begins by extracting these expert-defined rules from the active monitoring system (Figure 2), where domain experts author the conditions that trigger maintenance actions. Each rule  $\mathcal{R}^i$  typically comprises three key components: (1) a set of **conditions** that must be satisfied for a specified duration  $(t^i)$  to activate the rule. Let  $\mathcal{C}^i = \langle c_1^i, c_2^i, \dots, c_n^i \rangle$  be the set of atomic boolean conditions associated with it, where each  $c_j^i$  is a predicate over sensor readings or asset states (e.g., Temperature  $> 80^{\circ}\text{F}$ , Enthalpy  $< 15 \, \text{BTU/lb}$ ); (2) a set of **maintenance actions** for a rule  $(\mathcal{O}^i = \{o_1^i, o_2^i, \dots, o_{N_O^i}^i\})$  where  $N_O^i$  represents the total number of actions), hypothesized for verification once triggered; and (3) **metadata** including rule id  $(rule\_id)$ , rule name, asset type  $asset\_type$  (e.g., Air Compressor, Boiler),

From the domain documentation, we extract and assemble these conditions  $(\mathcal{C}^i)$  into a structured logical formula, referred to as a *condition tree*  $\mathcal{TR}^i$ . This tree is a boolean expression constructed from  $\mathcal{C}^i$  using logical operators  $\land$  (AND),  $\lor$  (OR), and optionally  $\neg$  (NOT). Formally: i) Each leaf node in  $\mathcal{TR}^i$  corresponds to an atomic condition  $c^i_j$ , ii) Internal nodes represent logical operators from the set  $\{\land,\lor,\lnot\}$ , iii) The root node evaluates to True if the entire condition tree is satisfied given the current sensor state. Given the above formulation, we define a rule  $\mathcal{R}^i$  as the tuple:

rule description, and estimated cost savings (in dollars) achieved by applying rule  $(\mathcal{C}^i)$ .

$$\mathcal{R}^i = (rule\_id, asset\_type, \mathcal{TR}^i, \mathcal{O}^i, \mathcal{C}^i, t^i)$$

We denote the collection of expert-written rules denoted as  $\mathcal{DS}_{\mathcal{R}} = \{\mathcal{R}^i\}_{i=1}^{N_{\mathcal{R}}}$ , where  $N_{\mathcal{R}} = 120$ . The rules span several asset types (10+) listed in Table 3. Table 3 summarizes, for each asset type, the number of rules (#Rules), the number of disjunction operators (# $\vee$ ), the total number of atomic conditions (# $\mathcal{C}$ ), and the number of observations. The count of disjunctions (# $\vee$ ) is particularly informative, as it reflects the branching complexity within the condition trees  $\mathcal{TR}^i$ , enabling us to sample a diverse range of conditions to which each rule  $\mathcal{R}^i$  applies.

To provide additional context on these industrial assets, we developed concise descriptions (Desc) for each asset type in collaboration with industry experts (Appendix Table 6). These descriptions are incorporated into the question-generation process to improve the relevance and clarity of the dataset.

#### 3.2 QA GENERATION PIPELINE

We design two primary types of diagnostic questions to evaluate a language model's reasoning capabilities under varying constraints. The first type requires the model to identify the **most relevant** option given the Question Conditions QC, testing its ability to prioritize the most probable root cause based on domain-specific knowledge. The second type requires selecting the **least relevant** option, challenging the model to distinguish between superficially similar answers and recognize when a condition-action mapping is unsupported by the available evidence. We refer to these question types as **selection** and **elimination**, respectively. Prior to question generation, we compute the following metadata to support dataset construction:

- i) Rule-Rule Similarity (RRSim): For each rule  $\mathcal{R}^i$ , we construct a textual representation of its condition tree  $\mathcal{TR}^i$  based on expert-authored documentation. We then generate embeddings for these texts and calculate pairwise cosine similarity scores across all rules. This metric enables sampling of semantically similar or diverse rules when constructing questions.
- ii) Unique Observations (UO): We manually curate and categorize observations/actions across the rule set  $\mathcal{DS}_{\mathcal{R}}$  to identify a universal set of unique observations. These form a candidate pool for selecting answer options. At present, |UO| = 193 and average length of each observation is around 20 (See Appendix 10 for distribution).

#### 3.2.1 Question-Answer Structure

Each question  $Q^i$  in the dataset is represented as a tuple (AD,QC,QP,OPT,A), where AD denotes the asset name along with its description obtained from Desc, QC specifies the observed conditions exhibited by the asset in the context of the question,QP represents the question prompt; further details on its construction are provided below, OPT is the set of answer options around 4 or 10. and A indicates the ground-truth correct answer for the question (Single true in at present).

#### 3.2.2 Rule Representation to Disjunctive Normal Form (DNF)

The QA generation procedure is summarized in Algorithm 1 (Appendix K). For each rule  $\mathcal{R}^i$ , we first convert the condition tree  $\mathcal{TR}^i$  into its *Disjunctive Normal Form (DNF)* that is, a disjunction (OR) over conjunctions (ANDs) of atomic conditions:

$$\mathcal{TR}_{ ext{DNF}}^i = igvee_{k=1}^K \left(igwedge_{j=1}^{m_k} c_{kj}^i
ight)$$

Each conjunctive clause in this DNF represents a *complete and specific observation pattern* sufficient to activate the rule  $\mathcal{R}^i$ . This formulation allows us to consider K distinct conjunctions that satisfy  $\mathcal{TR}^i_{\mathrm{DNF}} = \mathrm{True}$ .

For example, consider a rule  $\mathcal{R}^0$  with atomic conditions:  $c_1 = (\text{Preheat Valve}\% \geq 5\%), c_2 = (\text{Heating Valve}\% \geq 5\%)$  and  $c_3 = (\text{Heating Drained Flag} = 1 \text{ if reporting})$  let the condition tree be  $\mathcal{TR}^0 = (c_1 \vee c_2) \wedge c_3$  and the and its DNF form  $(c_1 \wedge c_3) \vee (c_2 \wedge c_3)$ . Each conjunction,  $(c_1 \wedge c_3)$  and  $(c_2 \wedge c_3)$ , is treated as a distinct, fully instantiated observation scenario, which is used as a unique QC for generating a question. This systematic transformation ensures that each QA instance is grounded in logically valid asset state combinations, reflecting domain expert reasoning and enabling scalable, interpretable dataset construction.

#### 3.2.3 SELECTING SETS OF ACTIONS

Next, we select observation combinations to construct the answer options for both *selection* and *elimination* question types. For *selection*-type questions (extracted\_obs\_sel in Algorithm 1), we identify candidate incorrect options by retrieving the  $N_{sel\_topk}$  least similar rules to  $\mathcal{R}^i$  using **RRSim** and collecting their observations, denoted as  $INC^i = \{inc^i_j\}_{j=1}^{N_{inc}}$ . We then generate answer tuples  $\{(QP^i_j, OPT^i_j, A^i_j)\}_{j=1}^{N_{sel}}$  by:

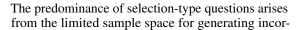
- (1) Selecting each  $o_i^i \in \mathcal{O}^i$  as the correct option, (2) Random sample  $\alpha$  incorrect options  $\in INC^i \setminus \mathcal{O}^i$ ,
- (3) Composing the prompt  $QP_i^i$ , which is drawn randomly from a pool of  $N_{QT}$  question templates.

For elimination-type questions (extracted\_obs\_eli in Algorithm 1), the correct options correspond to observations that do not belong to  $\mathcal{R}^i$ . Specifically, we retrieve the  $N_{ele\_topk}$  least similar rules to  $\mathcal{R}^i$  using **RRSim** and collect their observations, denoted as  $COR^i = \{cor_j^i\}_{j=1}^{N_{cor}}$ . We then randomly sample  $\min(N_{cor},\beta)$  observations from  $COR^i$  as correct options. For each correct option, we construct a question by pairing it with incorrect options sampled from  $\mathcal{O}^i$ , ensuring the elimination task challenges the model to identify irrelevant actions in the context of the given conditions.

 $\alpha$  and  $\beta$  are hyperparameters controlling the number of questions per rule. Larger values increase question count but reduce diversity, while smaller values enhance uniqueness but limit coverage.

#### 4 DIAGNOSTICIQ

We apply the pipeline described in 3 on 120 expert curated rules to create DiagnosticIQ. we set the hyperparameters  $N_{sel.topk}=25$ ,  $N_{eli.topk}=25$ ,  $N_{QT}=10$ ,  $\alpha=4$  and  $\beta=4$  during the creation process. The final dataset contains 6690 questions, with asset composition detailed in Figure 3. Selection-based questions make up 77.4% of the dataset, compared to 22.6% elimination questions. As shown in Figure 3 the majority of QA instances focus on AHU related scenarios (58.2%), followed by Chiller and Fan (5.9%)



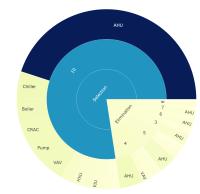


Figure 3: DiagnosticIQPro Composition by Asset/Number of Options/Question Type

rect options in elimination questions, which rely on the set  $\mathcal{O}^i$  for a given  $\mathcal{R}^i$ . Although this imbalance could be adjusted by setting  $\alpha \ll \beta$ , as previously discussed, doing so risks generating many similar elimination questions, ultimately reducing the diversity of the dataset.

The dataset intentionally reflects a real-world class imbalance, with a majority of rules pertaining to Air Handling Units (AHUs), as shown in Figure 3. We preserve this skew rather than resampling to accurately model operational priorities. Furthermore, while 'selection' questions have a fixed size of 10 options, 'elimination' questions feature a variable number of choices. This is because the incorrect options for elimination questions are drawn from the smaller, contextually relevant pool of alternatives available within the original source rule.

```
1 {
2    "AD": "Closed-loop water-cooled chiller system with cooling tower.",
3    "QC": ["Chiller Running", "Evaporator Delta T < 7degF",
4    "Cooling Tower Supply Temp < Setpoint - 4degF", "OAT > 43degF"],
5    "QP": "Given the observed conditions, what is the most likely root cause?",
6    "OPT": ["(A) Cooling tower is overcooling condenser water",
7    "(B) Chiller evaporator is fouled",
8    "(C) Supply water pump is cavitating",
9    "(D) Building load is too low"],
10    "A": "(A)"
```

Listing 1: DiagnosisIQ QA Instance for Chiller-Tower Case

#### 4.1 VARIANTS

We tweak DiagnosticIQdataset to create several variants that test LLMs under different conditions.

- (1) **DiagnosticIQPro**: We features a wider range of answer choices, with 10-option (increasing  $\alpha$  and  $\beta$  up to 10) questions constituting the majority (77.4%), thereby increasing dataset complexity (Figure 3). This distribution aligns with both practical asset relevance and intentional design decisions to balance diagnostic depth with question difficulty.
- (2) **DiagnosticIQPert**: We Perturb MCQ in DiagnosticIQ utilizing PertEval benchmark Li et al. (2024a) (Appendix K.3). We use this dataset to evaluate robustness against formatting.
- (3) **DiagnosticIQRationale**: We develop a dataset that has the rationale the LLM follow to arrive at the answer for MCQA questions, this is used to conduct human evaluation of the expert knowledge an LLM posses in the domain of maintenance recommendation. Further we utilize this for finetuning for the generalizability study.
- (4) **DiagnosticIQVerbose**: To identify the effect of presenting QC in natural language, we develop a variant of DiagnosticIQ that converts the symbolic representation of QC to natural language. The procedure for conversion can be found in (Appendix K.4)

#### 5 EXPERIMENTAL RESULTS

We perform a direct zero-shot prompting of the generated questions to assess the reasoning capacity of the LLMs (representative examples are provided in Appendix Figure 9). Our evaluation identifies areas of underperformance under zero-shot conditions. For evaluation we consider the Accuracy and the Macro-Accuracy as the main evaluation metrics. We utilize the Macro accuracy as DiagnosticIQ has a Asset class imbalance and will be used. Accuracy as  $\operatorname{Acc} = \frac{1}{|D_{\mathcal{Q}}|} \sum_{x \in D_{\mathcal{Q}}} \mathbb{1}[M(q(x)) = y_x]$  and Macro accuracy as  $\operatorname{Acc}_{\mathrm{macro}} = \frac{1}{|A|} \sum_{a \in A} \frac{1}{|D_a|} \sum_{x \in D_a} \mathbb{1}[M(q(x)) = y_x]$  where q(x) is the generated prompt, M(q(x)) the model's response,  $\mathbb{1}[\cdot]$  the indicator function returning 1 for correct predictions and 0 otherwise,  $D_{\mathcal{Q}}$  represents the whole dataset, A is the set of assets and  $D_a$  represents the questions belonging to asset class a.

#### 5.1 LEARDERBOARD

The comparative results in Table 1 show Claude-3-7-Sonnet achieving the highest Macro accuracy on both tasks (70.61% on DiagnosisIQ and 56.63% on DiagnosisIQPro), with larger models like Mistral-Large also performing well, while smaller models such as LlaMA-3-1-8B lag behind, indicating a

Table 1: Leaderboard: DiagnosisIQ and Pro.

Rank	Model	Macro. Diag IQ	Macro. +Pro	Diag IQ	+Pro
1	claude-3-7-sonnet*	70.61	56.63	72.66	53.80
2	deepseek-v3	67.02	41.38	67.89	35.80
3	o1*	65.41	24.79	70.22	26.11
4	mistral-large	63.15	41.13	65.52	36.50
5	qwen2-5-72b-ins.	61.22	35.91	63.09	32.93
6	llama-3-3-70b-ins.	61.67	36.56	60.33	32.27
7	mistral-small-3-1-24b	61.17	33.79	60.15	28.42
8	mistral-medium-2505	60.34	35.36	61.43	30.16
9	granite-3-3-8b-instruct	59.45	42.39	57.26	31.43
10	gemini-2.0-flash*	57.64	26.65	54.63	20.82
11	llama-3-1-405b-ins.	56.56	38.82	59.03	35.58
12	gemini-1.5-pro*	53.14	24.72	65.44	27.77
13	microsoft-phi-4	50.52	31.35	47.50	23.99
14	claude-3-5-haiku*	46.93	17.72	44.41	15.55
15	llama-3-1-8b-ins.	38.69	18.80	36.70	12.89

clear correlation between model size and performance. These findings establish a strong baseline, revealing that general-purpose LLMs struggle with reasoning about sensor conditions in industrial settings, and the sharp drop in accuracy on DiagnosisIQPro highlights the challenge of larger, realistic action spaces. Overall, the leaderboard demonstrates a pronounced performance gap favoring larger models and underscores the critical need for domain-specific knowledge integration to enable effective real-world industrial diagnostics.

Apart for Claude-3-7-Sonnet, the performance of all remaining 10 models on the DiagnosisIQPro dataset is below 45%, revealing a surprising gap on this complex task and underscoring the universal need for improvement in handling industrial diagnostic reasoning.

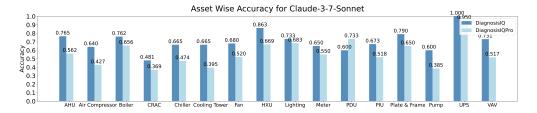


Figure 4: The Accuracy variation across varies industrial assets for claude-3-7-sonnet

Asset-wise Analysis We select Claude-3-7-sonnet from the leaderboard for asset-wise analysis. Figure 4 shows asset-wise accuracy comparisons, revealing consistently higher accuracy on assets like UPS 100.0%) and HXU 86.3%) in DiagnosisIQ, but analysing DiagnosisIQPro considerably drops overall, where HXU drops by -19.4% (to 66.9%). Further for CRAC (-11.2%), Cooling Tower (-27.0%), and Pump (-21.5%). These results highlight the model's domain sensitivity, with some assets maintaining robust performance, and expose a reasoning complexity gap, where accuracy declines sharply in more challenging, multi-condition scenarios.

Question type wise Analysis. Figure 12 compares model accuracy on selection and elimination tasks across both datasets. Elimination consistently yields higher accuracy (e.g., Mistral-Large achieves 71.0% on elimination-DiagnosisIQ vs. 63.9% on selection-DiagnosisIQ). However, the performance drop from DiagnosisIQ to DiagnosisIQPro is more pronounced for selection questions (Mistral-Large drops from 63.9% to 28.9%) than for elimination questions (71.0% to 62.6%), indicating that compositional reasoning in complex scenarios is particularly challenging when models must select the most relevant option. The same patterns seems to be consistently shown in multiple models.

**Robustness Against Perturbation.** MCQA datasets inherently contain bias due to planted correct answers. To assess this, we applied perturbation analysis on DiagnosisIQ using the PertEval benchmark Li et al. (2024a) (Appendix K.3), generating perturbed questions and measuring

379

380

381

382

384 385

386 387

389 390 391

392 393

394

396 397

398

399

400

401

402

403

404

405

406

407

408 409

410

411

412

413

414

415 416 417

418

419

420

421

422

423

424

425

426

427

428 429

430

431

accuracy (Acc) on both original and perturbed sets. Consistency accuracy was computed as  $Acc@Consist = \frac{1}{|D_{\mathcal{Q}}|} \sum_x C(x) \wedge C_{\text{perturb}}(x)$ , where  $C_{\text{perturb}}(x)$  indicates correct predictions on perturbed prompts. Table 2 reports Acc@Perturb, perturbation drop rate (PDR), and Acc@Consist, highlighting variation in model robustness to question perturbations.

#### 5.2 EVALUATING DOMAIN UNDERSTANDING AND REASONING

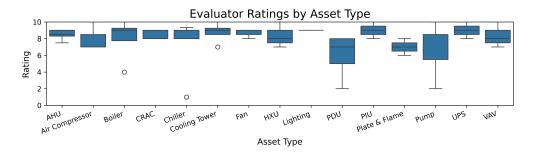


Figure 5: The expert rating for reasoning under mistral-large

To assess the domain understanding and reasoning patterns of LLMs, we conducted a human evaluation of the rationales they generate for maintenance recommendations. We prompted Mistral-Large to provide its reasoning given a question and correct answer on a representative sample of 27 questions from DiagnosticIQ, ensuring at least one question from each asset type. We explicitly provide the correct answer as we are specifically evaluating whether the reasoning patterns of an LLM can match a domain expert. We then had five domain

Table 2: PertEval, Significant \*\* ( $\alpha = 0.01$ ) Model Acc@Per. PDR Acc@Con. llama-3-3-70b 66.77 0.11\*\* 52.30 -0.02\*\* deepseek-v3 66.70 57.66 gwen2-5-72b 62.61 -0.01 52.28 llama-3-1-405b 0.03\*\* 60.66 48.09 -0.02\*\* mistral-medium 60.00 49.45 -0.12\*\* mistral-large 57.39 49.92 -0.04\*\* micro.-phi-4 45.36 32.00 llama-3-1-8b 44.03 0.20\*\* 23.01

experts rate each generated rationale on a scale of 0 (incorrect) to 10 (expert-level quality).

As shown in Figure 5, the analysis reveals that the model's reasoning generally aligns with expert reasoning, supporting its potential to augment maintenance tasks. However, for certain assets (e.g., PDU, Pump, Boiler), we observe inter-rater disagreement, with expert feedback indicating differing expectations on the required granularity of the explanation. Furthermore, the model consistently scored lower among evaluators on the Plate & Frame HX asset, with experts expressing that the rationales lack the nuanced operational knowledge required, signifying a potential knowledge gap for that specific equipment type. These concerns must be addressed when deploying LLMs in this task.

#### 5.3 GENERALIZABILITY STUDY

Transfer learning between different shows promising results industrial automation Maschler & Weyrich (2021) for tasks such as fault prediction or anomaly detection. However, transfer learning of rules between different industrial assets is an interesting direction that has not been studied. We consider Qwen3-8B Yang et al. (2025), Llama-3.1-8B-Instruct Dubey et al. (2024), and Granite-3.3-8B-Instruct Granite Team (2024). To avoid any data leakage, the split is stratified by asset. To account for the imbalance in the questions per asset (Figure 3), we consider two splits: AHU and the rest of the assets. For each model we finetune on each split and test on the other split. We use Supervised Fine-Tuning (SFT) and GRPO. Overall, the rule learning is transferable across assets, with Qwen3-8B being the best performing model both with and without SFT. More details on the prompt formatting in Appendix Section F. Table 6 provides the results.

#### 5.4 CONDITION FORMATTING STUDY

We investigate effect converting the temporal conditions into natural language as a pre-processing step as in Zhang et al. (2025). We utilize the Macro accuracies of the DiagnosticIQVerbose dataset

Figure 6: SFT/GRPO experiments on training/testing AHII/Other splits

ing/testing Arro/Other spints				
Model	Setting	AHU	Other	
Llama-3.1-8B	Base	50.88	44.52	
	SFT	52.31	56.44	
	GRPO	52.95	61.76	
Qwen3-8B	Base	56.47	66.85	
	SFT	68.89	80.28	
	GRPO	55.27	64.76	
Granite-3.3-8B	Base	59.02	56.58	
	SFT	58.79	59.61	
	GRPO	54.94	63.76	

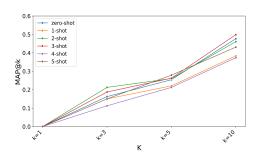


Figure 7: MAReE MAP@K with number of examples from DiagnosticIQ

(dataset creation details can be found in Appendix K.4) and the difference in Macro accuracy of DiagnosticIQ and DiagnosticIQVerbose ( $\Delta$  Macro Accuracy) for this and compare the effectiveness of the formatting. We present our findings in 4. We identify that the Macro accuracies drop on almost all models which shows the symbolic understanding the LLMs possess atleast in the domain of physical asset maintenance.

#### 5.5 Maintenance Action Rec. Engine

We present Maintainance Action Recommendation Engine (MAReE), a deployed application leveraging DiagnosticIQ to recommend maintenance actions based on abnormal conditions defined by subject matter experts (SMEs). MAReE employs LLM-Score, which assigns relevance scores to candidate actions. We evaluate MAReE on 11 new SME-authored rules with ground truth actions, varying  $k \in \{1, 3, 5, 10\}$  and measuring MAP@K (Mean Average Precision at k), indicating whether the ground truth action appears in the top-k recommendations. For each rule, 10 candidate actions are dynamically sampled using an embedding model to create realistic and challenging evaluation scenarios. Results in Figure 7 show that 3-shot prompting achieves the highest MAP@10 score (49.84%), outperforming zero-shot and other shot counts, while MAP@K scores for K > 1 fluctuate without a clear trend, suggesting that increasing example count does not consistently improve recall beyond the top prediction, though few-shots beat the zero-shot at any k. Detailed system prompt is available in Appendix 11.

#### 6 CONCLUSION AND LIMITATIONS

This paper addresses the gap of systematic evaluation of LLMs on maintenance action recommendation in the industrial setting. We develop a generic pipeline that systematically converts symbolic, human-authored rules into a MCQA format. Utilizing our pipeline we introduce DiagnosticIQ and its variants designed to benchmark the ability of utilizing LLMs to recommend maintenance actions. We benchmark 15 leading LLMs, establishing the first standardized leaderboard for this task. Our analysis systematically evaluated model reasoning, generalization, and robustness, providing a clear baseline for the community.

Our work confirms the potential of LLMs as a powerful tool in this domain and provide resources for further developing reliable industrial AI. However, we acknowledge a few key limitations and future directions of research. The industrial domain includes over 800 asset types, but our current study covers only a limited subset again due to the extensive amount of time being spent (nearly 120 days using 3-4 SMEs) on writing each pair of condition-rule and action. We plan to expand assets coverage to improve robustness and mitigate potential biases due to LLM familiarity with specific assets. Furthermore, although out dataset generation pipeline is generic and applicable to other domains such as business process management and cloud resource monitoring our experiments so far focus on Industry 4.0. Extending this approach to additional domains remains an important future direction. We envision this work as foundational step towards more sophisticated assistive tools that augment rule creation as a whole and automated asset monitoring. The MAReE experiments underscores the need for more advance methods or finetuned models to reliably address this task

#### REFERENCES

- Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang Li, Sheila Babayan, Kavya Kopparapu, Zachary Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srinivasan, Manzil Zaheer, Felix Yu, and Sanjiv Kumar. Rest meets react: Self-improvement for multi-step reasoning llm agent, 2023. URL https://arxiv.org/abs/2312.10003.
- A Beghi, R Brignoli, Luca Cecchinato, Gabriele Menegazzo, Mirco Rampazzo, and F Simmini. Data-driven fault detection and diagnosis for hvac water chillers. *Control Engineering Practice*, 53:79–91, 2016.
- Sathvik Bhaskarpandit. Wind power forecasting, 2020. URL https://www.kaggle.com/datasets/theforcecoder/wind-power-forecasting/data.
- Bishwaranjan Bhattacharjee, Aashka Trivedi, Masayasu Muraoka, Muthukumaran Ramasubramanian, Takuma Udagawa, Iksha Gurung, Nishan Pantha, Rong Zhang, Bharath Dandala, Rahul Ramachandran, Manil Maskey, Kaylin Bugbee, Michael M. Little, Elizabeth Fancher, Irina Gerasimov, Armin Mehrabian, Lauren Sanders, Sylvain V. Costes, Sergi Blanco-Cuaresma, Kelly Lockhart, Thomas Allen, Felix Grezes, Megan Ansdell, Alberto Accomazzi, Yousef El-Kurdi, Davis Wertheimer, Birgit Pfitzmann, Cesar Berrospi Ramis, Michele Dolfi, Rafael Teixeira De Lima, Panagiotis Vagenas, S. Karthik Mukkavilli, Peter W. J. Staar, Sanaz Vahidinia, Ryan McGranaghan, and Tsengdar J. Lee. INDUS: Effective and efficient language models for scientific applications. In Franck Dernoncourt, Daniel Preoţiuc-Pietro, and Anastasia Shimorina (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 98–112, Miami, Florida, US, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-industry.9. URL https://aclanthology.org/2024.emnlp-industry.9/.
- Jian Chen, Peilin Zhou, Yining Hua, Loh Xin, Kehui Chen, Ziyuan Li, Bing Zhu, and Junwei Liang. FinTextQA: A dataset for long-form financial question answering. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6025–6047, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.328. URL https://aclanthology.org/2024.acl-long.328/.
- Christodoulos Constantinides, Dhaval Patel, Shuxin Lin, Claudio Guerrero, Sunil Dagajirao Patil, and Jayant Kalagnanam. Failuresensoriq: A multi-choice qa dataset for understanding sensor relationships and failure modes, 2025. URL https://arxiv.org/abs/2506.03278.
- OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models. https://github.com/open-compass/opencompass, 2023.
- Andrew A Cook, Göksel Mısırlı, and Zhong Fan. Anomaly detection for iot time-series data: A survey. *IEEE Internet of Things Journal*, 7(7):6481–6494, 2019.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms, 2023. URL https://arxiv.org/abs/2305.14314.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Shanghua Gao. CURE-Bench-InternalReasoning. https://kaggle.com/competitions/cure-bench-internal-reasoning, 2025. Kaggle.
- IBM Granite Team. Granite 3.0 language models, October 2024. URL https://github.com/ibm-granite/granite-3.0-language-models/.
- Sungwon Han, Jinsung Yoon, Sercan O Arik, and Tomas Pfister. Large language models can automatically engineer features for few-shot tabular learning. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 17454–17479. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/han24f.html.

Seongtae Hong, Joong Min Shin, Jaehyung Seo, Taemin Lee, Jeongbae Park, Cho Man Young, Byeongho Choi, and Heuiseok Lim. Intelligent predictive maintenance RAG framework for power plants: Enhancing QA with StyleDFS and domain specific instruction tuning. In Franck Dernoncourt, Daniel Preoţiuc-Pietro, and Anastasia Shimorina (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 805–820, Miami, Florida, US, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-industry.61. URL https://aclanthology.org/2024.emnlp-industry.61/.

- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
- Xi Jiang, Jian Li, Hanqiu Deng, Yong Liu, Bin-Bin Gao, Yifeng Zhou, Jialin Li, Chengjie Wang, and Feng Zheng. MMAD: A comprehensive benchmark for multimodal large language models in industrial anomaly detection. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=JDiER86r8v.
- Ameeth Kanawaday and Aditya Sane. Machine learning for predictive maintenance of industrial machines using iot sensor data. In 2017 8th IEEE international conference on software engineering and service science (ICSESS), pp. 87–90. IEEE, 2017.
- Sunwoo Lee, Dhammiko Arya, Seung-Mo Cho, Gyoung-eun Han, Seokyoung Hong, Wonbeom Jang, Seojin Lee, Sohee Park, Sereimony Sek, Injee Song, Sungbin Yoon, and Eric Davis. Tel-Bench: A benchmark for evaluating telco-specific large language models. In Franck Dernoncourt, Daniel Preoţiuc-Pietro, and Anastasia Shimorina (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 609–626, Miami, Florida, US, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. emnlp-industry.45. URL https://aclanthology.org/2024.emnlp-industry.45/.
- Jiatong Li, Renjun Hu, Kunzhe Huang, Yan Zhuang, Qi Liu, Mengxiao Zhu, Xing Shi, and Wei Lin. Perteval: Unveiling real knowledge capacity of llms with knowledge-invariant perturbations, 2024a. URL https://arxiv.org/abs/2405.19740.
- Ruosen Li, Zimu Wang, Son Quoc Tran, Lei Xia, and Xinya Du. MEQA: A benchmark for multi-hop event-centric question answering with explanations. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024b. URL https://openreview.net/forum?id=snNuvAOQxB.
- Yahan Li, Keith Harrigian, Ayah Zirikly, and Mark Dredze. Are clinical t5 models better for clinical text?, 2024c. URL https://arxiv.org/abs/2412.05845.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2022. URL https://arxiv.org/abs/2109.07958.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025.
- Benjamin Maschler and Michael Weyrich. Deep transfer learning for industrial automation: A review and discussion of new techniques for data-driven machine learning. *IEEE Industrial Electronics Magazine*, 15(2):65–75, 2021.
- Alexander Nikitin and Samuel Kaski. Human-in-the-loop large-scale predictive maintenance of workstations. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3682–3690, 2022.
- Ahmed Okudan. Predictive maintenance dataset air compressor, 2023. URL https://www.kaggle.com/datasets/afumetto/predictive-maintenance-dataset-air-compressor/data.
- OpenAI. o1-preview, 2024. URL https://openai.com/.

- Oracle. Add failure diagnostics information to asset incidents and anomalies. https://docs.oracle.com/en/cloud/saas/iot-asset-cloud/iotaa/add-failure-diagnostics-information-asset-incidents-and-anomalies.html, 2025. Oracle IoT Asset Monitoring Cloud Service.
- Sreshta R Putchala, Rithik Kotha, Vanitha Guda, and Yellasiri Ramadevi. Transformer data analysis for predictive maintenance. In *Proceedings of Second International Conference on Advances in Computer Engineering and Communication Systems: ICACECS 2021*, pp. 217–230. Springer, 2022.
- Sourjyadip Ray, Kushal Gupta, Soumi Kundu, Dr Payal Arvind Kasat, Somak Aditya, and Pawan Goyal. ERVQA: A dataset to benchmark the readiness of large vision language models in hospital environments. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 15594–15608, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.873. URL https://aclanthology.org/2024.emnlp-main.873/.
- M. Raza, Z. Jahangir, M. B. Riaz, et al. Industrial applications of large language models. *Scientific Reports*, 15:13755, 2025. doi: 10.1038/s41598-025-98483-1. URL https://doi.org/10.1038/s41598-025-98483-1.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof qa benchmark, 2023. URL https://arxiv.org/abs/2311.12022.
- Joshua Robinson, Christopher Michael Rytting, and David Wingate. Leveraging large language models for multiple choice question answering, 2023. URL https://arxiv.org/abs/2210.12353.
- Tobias Schimanski, Jingwei Ni, Roberto Spacey Martín, Nicola Ranger, and Markus Leippold. ClimRetrieve: A benchmarking dataset for information retrieval from corporate climate disclosures. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 17509–17524, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.969. URL https://aclanthology.org/2024.emnlp-main.969/.
- Gauri Shah and Aashis Tiwari. Anomaly detection in iiot: A case study using machine learning. In *Proceedings of the ACM India joint international conference on data science and management of data*, pp. 295–300, 2018.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv* preprint arXiv:2402.03300, 2024.
- Michael Stewart, Melinda Hodkiewicz, and Sirui Li. Large language models for failure mode classification: An investigation, 2023. URL https://arxiv.org/abs/2309.08181.
- Jing Su, Chufeng Jiang, Xin Jin, Yuxin Qiao, Tingsong Xiao, Hongda Ma, Rong Wei, Zhi Jing, Jiajun Xu, and Junhong Lin. Large language models for forecasting and anomaly detection: A systematic literature review, 2024. URL https://arxiv.org/abs/2402.10350.
- Ekaterina Sviridova, Anar Yeginbergen, Ainara Estarrona, Elena Cabrio, Serena Villata, and Rodrigo Agerri. CasiMedicos-arg: A medical question answering dataset annotated with explanatory argumentative structures. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 18463–18475, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1026. URL https://aclanthology.org/2024.emnlp-main.1026/.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, and etc. Gemini: A family of highly capable multimodal models, 2024. URL https://arxiv.org/abs/2312.11805.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL https://arxiv.org/abs/2307.09288.

- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models, 2023. URL https://arxiv.org/abs/2305.04091.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. MMLU-pro: A more robust and challenging multi-task language understanding benchmark. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=y10DM6R2r3.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese, John Schulman, and William Fedus. Measuring short-form factuality in large language models, 2024. URL https://arxiv.org/abs/2411.04368.
- Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuang Li, Binxin Hu, Wendy Gao, Jiaxing Xu, Yiming Liu, Jie Tang, Hongning Wang, and Minlie Huang. Benchmarking complex instruction-following with multiple constraints composition. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=U2aVNDrZGx.
- Ran Xu, Wenqi Shi, Yue Yu, Yuchen Zhuang, Yanqiao Zhu, May Dongmei Wang, Joyce C. Ho, Chao Zhang, and Carl Yang. BMRetriever: Tuning large language models as better biomedical text retrievers. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 22234–22254, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1241. URL https://aclanthology.org/2024.emnlp-main.1241/.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Hong Yang, Aidan LaBella, and Travis Desell. Predictive maintenance for general aviation using convolutional transformers, 2022. URL https://arxiv.org/abs/2110.03757.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *ArXiv*, abs/2210.03629, 2022. URL https://api.semanticscholar.org/CorpusID:252762395.
- Fanghua Ye, Mingming Yang, Jianhui Pang, Longyue Wang, Derek Wong, Emine Yilmaz, Shuming Shi, and Zhaopeng Tu. Benchmarking llms via uncertainty quantification. *Advances in Neural Information Processing Systems*, 37:15356–15385, 2024.

Jiahao Ying, Yixin Cao, Yushi Bai, Qianru Sun, Bo Wang, Wei Tang, Zhaojun Ding, Yizhe Yang, Xuanjing Huang, and Shuicheng YAN. Automating dataset updates towards reliable and timely evaluation of large language models. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=EvEqYlQv8T.

Tianyang Zhang, Zhuoxuan Jiang, Shengguang Bai, Tianrui Zhang, Lin Lin, Yang Liu, and Jiawei Ren. RAG4ITOps: A supervised fine-tunable and comprehensive RAG framework for IT operations and maintenance. In Franck Dernoncourt, Daniel Preoţiuc-Pietro, and Anastasia Shimorina (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pp. 738–754, Miami, Florida, US, November 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-industry.56. URL https://aclanthology.org/2024.emnlp-industry.56/.

Wen Zhang, Long Jin, Yushan Zhu, Jiaoyan Chen, Zhiwei Huang, Junjie Wang, Yin Hua, Lei Liang, and Huajun Chen. Trustuqa: A trustful framework for unified structured data question answering, 2024b. URL https://arxiv.org/abs/2406.18916.

Yudi Zhang, Pei Xiao, Lu Wang, Chaoyun Zhang, Meng Fang, Yali Du, Yevgeniy Puzyrev, Randolph Yao, Si Qin, Qingwei Lin, Mykola Pechenizkiy, Dongmei Zhang, Saravan Rajmohan, and Qi Zhang. RuAG: Learned-rule-augmented generation for large language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=BpIbnXWfhL.

Ziyin Zhang, Zhaokun Jiang, Lizhen Xu, Hongkun Hao, and Rui Wang. Multiple-choice questions are efficient and robust llm evaluators, 2024c. URL https://arxiv.org/abs/2405.11966.

Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are not robust multiple choice selectors, 2024. URL https://arxiv.org/abs/2309.03882.

Keyi Zhong, Tom Jackson, Andrew West, and Georgina Cosma. Natural language processing approaches in industrial maintenance: A systematic literature review. *Procedia Computer Science*, 232:2082–2097, 2024. ISSN 1877-0509. doi: https://doi.org/10.1016/j.procs.2024.02.029. URL https://www.sciencedirect.com/science/article/pii/S1877050924002060. 5th International Conference on Industry 4.0 and Smart Manufacturing (ISM 2023).

Yizhang Zhu, Shiyin Du, Boyan Li, Yuyu Luo, and Nan Tang. Are large language models good statisticians? In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024. URL http://papers.nips.cc/paper\_files/paper/2024/hash/729786203d330da046dd8091c2d92a66-Abstract-Datasets\_and\_Benchmarks\_Track.html.

## A DECLARATION OF GENERATIVE AI AND AI-ASSISTED TECHNOLOGIES IN THE WRITING PROCESS

During the preparation of this work, the authors used Grammarly in order to improve the grammar, clarity, and flow of the manuscript. After using this tool/service, the author(s) reviewed and edited the content as needed and take full responsibility for the content of the publication.

### B RULE GENERATION

The rules originally came from Smarter Buildings, which was introduced in 2011 as part of the larger Smarter Planet initiative. There were originally 18 Fault Detection and Diagnostic (FDD) rules spread across 3 equipment types, with Air Handlers being the primary focus. There were two main objectives to the program: to achieve 5-15% energy savings for the monitored equipment, and to reduce maintenance hours by 30%. Air Handlers were originally, and continue to be, the primary focus of the rules, as they are the most prevalent piece of equipment at any location and thus provide the greatest savings. As the program continued and the rule set expanded, additional equipment (such as Chillers and Boilers) was added for increased monitoring and savings, leading to a total of 118 active rules across 13 equipment types.

The rules were actively developed over the course of 7 years, with 11 major updates. Each update added additional rules and updated existing ones to account for updates to the logic. There were two key roles involved in developing the rules: the Reliability Engineer who developed the rules, and the System Administrator who coded them. There was also a Rules Logic committee that typically involved 2-8 participants that met every other week to brainstorm and develop the rough logic for the rules, with the Reliability Engineer and System Administrator working closely together to ensure the logic was coded correctly.

Writing a new rule does not take a long time to code – typically around 30 minutes to ensure it's running correctly. However, with additional testing and verification, it can take significantly longer to finalize. This includes documentation of the new rule, as well as the correct verbiage on the rule, to better inform technicians about the potential causes and where to begin troubleshooting. Updating a rule takes less time but still requires additional testing and updating of the documentation, so it is still not an insignificant task.

#### C LLMs in Industry 4.0

Large language models (LLMs) have seen rapid development and broad application across domains, from general-purpose models like OpenAI's GPT series (OpenAI, 2024) and Meta's Llama 2 (Touvron et al., 2023) to specialized, multimodal models such as Gemini (Team et al., 2024) and Mistral 7B (Jiang et al., 2023). These foundational models demonstrate impressive capabilities in natural language understanding, generation, and reasoning (Wei et al., 2022; Wang et al., 2023; Yao et al., 2022), and are increasingly benchmarked on complex question answering and reasoning datasets (Rein et al., 2023; Li et al., 2024b; Wang et al., 2024).

In the context of domain-specific applications, several efforts highlight the benefits of fine-tuning or training LLMs on domain-relevant corpora. For example, INDUS (Bhattacharjee et al., 2024) and TelBench (Lee et al., 2024) demonstrate improved performance on scientific and telecommunications tasks, respectively, underscoring the value of specialized data and benchmarks. Similarly, clinical text models (Li et al., 2024c) and biomedical retrievers (Xu et al., 2024) leverage domain-specific adaptations to better meet task requirements.

Industrial and predictive maintenance applications have attracted increasing attention with approaches combining LLMs and machine learning for failure mode classification and condition monitoring (Stewart et al., 2023; Nikitin & Kaski, 2022; Putchala et al., 2022; Yang et al., 2022). These studies highlight the potential of language models to extract actionable insights from maintenance logs, sensor data, and domain knowledge, aiding decision-making in complex industrial environments. Public datasets such as predictive maintenance for air compressors (Okudan, 2023) and wind power forecasting (Bhaskarpandit, 2020) facilitate research in this area.

Recent work also explores the robustness and reliability of LLMs in handling structured data, multihop reasoning, and multiple-choice question answering (MCQA). Studies have pointed out challenges in LLMs' MCQA performance (Robinson et al., 2023; Zhang et al., 2024c; Zheng et al., 2024) and proposed methods like chain-of-thought prompting (Wei et al., 2022) and plan-and-solve prompting (Wang et al., 2023) to elicit better reasoning. Self-improving multi-step reasoning agents (Aksitov et al., 2023) and trustful frameworks for unified question answering (Zhang et al., 2024b) further advance LLM capabilities.

Benchmarking platforms such as OpenCompass (Contributors, 2023) and uncertainty quantification in benchmarks (Ye et al., 2024) enable more reliable evaluation of LLMs across tasks and domains. Studies on knowledge capacity and factuality assessment (Li et al., 2024a; Wei et al., 2024) inform improvements in LLM trustworthiness.

In the tabular and scientific domain, LLMs have shown promise for automatic feature engineering (Han et al., 2024) and scientific knowledge extraction (Bhattacharjee et al., 2024). Efforts to automate dataset updates and maintain evaluation relevance (Ying et al., 2024) address practical challenges in large-scale LLM deployment.

#### D RULE DOCUMENT COLLECTED FROM INDUSTRY EXPERTS.

Table 3: Statistics of expert-curated rules collected (Asset type definitions can be found in Table 6).

Asset Type	#Rules	#\/	# <i>C</i>	#Observations
Fan	1	0	4	4
UPS	1	1	2	1
Lighting	1	1	2	2
Plate & Frame	1	1	4	4
PIU	2	1	5	5
Meter	3	4	4	4
Air Compressor	3	1	6	8
PDU	3	0	7	3
HXU	4	0	12	9
Cooling Tower	4	1	11	12
Pump	5	1	20	22
Boiler	6	5	17	19
VAV	8	3	30	18
CRAC	10	0	21	28
Chiller	11	0	26	27
AHU	55	54	312	172

#### E CONDITIONS AS NATURAL LANGUAGE EXPERIMENTS

Table 4: Conditions as Natural Language

Table 4. Collutions as Natural Language				
Model	Macro Accuracy	<b>△ Macro Accuracy</b>		
claude-3-7-sonnet	68.22	-2.58		
o1	64.53	-4.33		
deepseek-v3	62.02	-6.45		
mistral-large	60.45	-3.70		
qwen2-5-72b-instr.	56.73	-3.53		
llama-3-3-70b-instr.	55.16	-9.62		
mistral-medium	55.89	-5.09		
gemini-1.5-pro	55.78	-4.10		
mistral-small-3-1.	53.78	-6.01		
llama-3-1-405b-instr.	53.81	-4.37		
granite-3-3-8b-instr.	53.69	-1.75		
gemini-2.0-flash	49.78	-6.37		
microsoft-phi-4	45.64	-2.73		
claude-3-5-haiku	42.30	-6.22		
llama-3-1-8b-instr.	39.79	+1.69		

#### F GENERALIZABILITY EXPERIMENTS SETUP

**SFT:** For hardware we use 4xNvidia A100 GPUs with 80GB memory. We fine-tune the base model using QLoRA (Dettmers et al., 2023) with FlashAttention-2 (max sequence length 2048, packed sequences), 4-bit quantization, and LoRA adapters ( $r=16, \alpha=16$ ), training for 3 epochs with a per-device batch size of 8, a learning rate of  $2\times 10^{-4}$  and a 0.1 warmup ratio.

**GRPO:** For hardware we use 4xNvidia A100 GPUs with 80GB memory. We train for 250 steps, with 16 generations per step and effective batch size of 4 per device. We use Learning Rate (LR) of  $5*10^{-7}$ , Beta of 0.001, cosine LR scheduler and 0.03 warmup ratio. We use the Hugging Face implementation that excludes prompt length and reward std normalization due to bias Liu et al. (2025).

**Formatting:** For Qwen3-8B we align on a json format with reasoning and answer fields as recommended in the documentation for MCQA questions. For granite-3.3-8B we use <think></think> and <response></response> tags as described in the model card and in llama-3.1-8b we used <think></think> and <answer></answer> tags. Think tags/fields are ommitted for SFT. During evaluation, if a model doesn't provide an answer we consider it as wrong.

#### G BACKGROUND: INDUSTRIAL ASSET DIAGNOSTICS SYSTEM

Large corporations and institutions usually own their industrial facilities. Industrial facilities rely on a diverse set of physical assets, including but not limited to chillers, boilers, pumps, compressors, air compressors, and air handling units (AHUs), to maintain safe, efficient, and resilient operations and ensure a smooth working environment. These assets are foundational in sectors like data centers, hospitals, manufacturing plants, and commercial buildings, where equipment failures can lead to operational downtime, safety risks, or financial losses.

Today, to manage these industrial assets proactively, organizations deploy sensor networks that continuously track real-time measurements such as temperature, pressure, flow rate, humidity, and power consumption. Domain experts use this data to define *diagnostic rules*, which map specific combinations of sensor conditions to likely early signals of particular failure modes and ideally recommended follow-up actions, such as inspection, warning bookkeeping, and proactive maintenance. These rules power the early warning systems and suggest the predictive maintenance workflows, allowing operations teams to detect and address potential faults before they escalate to interrupt the operation.

Creating these diagnostic rules is a labor-intensive, highly specialized task. Domain experts must understand the behavior of each asset under various operating conditions, interpret complex sensor relationships, and encode domain knowledge into logical expressions that accurately represent the behavior of each asset. A single facility may require hundreds of rules per asset type or model, each reflecting detailed dependencies between multivariate sensor signals. Rules often involve temporal thresholds, conditional logic, and asset-specific tolerances.

Table 5 presents a set of representative rules derived from production environments, demonstrating the range of conditions and asset behaviors encountered in industrial diagnostics.

This rule-related management domain presents several challenges that complicate automation:

- High-dimensional, domain-specific sensor data, often with implicit semantics not found in general corpora.
- Complex logical dependencies across multiple sensor conditions, often involving nested boolean logic.
- Asset heterogeneity, where similar asset classes behave differently depending on configuration or environment.
- Tacit expert knowledge that is rarely documented and typically acquired through experience.

While diagnostic rules are effective in practice, they do not scale easily. The growing complexity and data richness of industrial systems require tools that can assist in generating, validating, and refining such rules

Large Language Models (LLMs) offer a promising avenue for this. However, general-purpose LLMs are not trained on sensor semantics or domain-specific diagnostics, and it is unclear whether they can reason over the kinds of multivariate conditions and logic used in real-world maintenance workflows.

Table 5: Illustrative Examples of Diagnostic and Alert Rules across Industrial Asset Types

Asset Type	Rule Name	Rule Logic Summary
AHU	Simultaneous Heating and Cooling	AHU Running; Cooling Valve ≥ 5%; Heating Valve ≥ 5%; Drain Flags Active; Met for 2 Hours
AHU	Heating Valve Open When Warm Outside	AHU Running; OAT - Supply Temp textmore 5°F; Heating Valve textmore 10%; Met for 2 Hours
Air Compressor	Pressure Setpoint Attainment	ABS(Pressure - Setpoint) textmore 10 PSI OR Pressure textmore 130 PSI (if setpoint missing); Met for 2 Hours
Air Compressor	Flow Flag	Not Monday; Air Flow textmore 120% of Previous Day's Average; Met for 2 Hours
Boiler	Excess O <sub>2</sub> in Stack	Fuel Flow textmore 5 and Flue Gas O <sub>2</sub> exceeds threshold; Met for 2 Hours
Boiler	Flue Gas Temperature Setpoint Attainment	Flue Gas Temp below setpoint; Met for 2 Hours
Chiller	Temperature Setpoint Attainment	Chiller Running; Supply Temp - Setpoint textmore 5°F; Met for 2 Hours
Chiller	Low Supply Temperature	Chiller Running; Setpoint - Supply Temp textmore 3°F; Met for 2 Hours
Cooling Tower	Delta T Out of Range	Condenser Return - Supply Temp <5°F; Tower Running; Met for 2 Hours
Cooling Tower	Pressure Setpoint Attainment	ABS(Condenser Pressure Diff - Setpoint) textmore 5 PSI; OAT <95°F; Met for 2 Hours

Given the domain complexity and the limitations of manual rule engineering, we now present Asset DiagnosticIQ, a benchmark for testing whether LLMs can assist in scalable, high-quality industrial diagnostics.

#### H ASSET TYPE DESCRIPTIONS

The asset DiagnosisIQ dataset includes diagnostic rules derived from a wide range of industrial asset types commonly found in commercial buildings, data centers, manufacturing facilities, and other operational environments. Each asset class is associated with domain-specific behaviors, sensor signals, and potential fault conditions that inform the construction of diagnostic questions. This section provides concise descriptions of the primary asset types represented in the dataset, supporting an understanding of their operational roles and diagnostic relevance.

Table 6 provides brief descriptions of the major asset types represented in the Asset DiagnosisIQ dataset. These physical systems are typically monitored via real-time sensor data and are subject to diagnostic rules used for fault detection and predictive maintenance.

# I CASE STUDY: CLOSED-LOOP WATER-COOLING CHILLER WITH COOLING TOWER

This case study illustrates a realistic multi-component diagnostic scenario for a **closed-loop water-cooled chiller system** paired with a **cooling tower**. It demonstrates how expert-authored rules, time-persistent sensor conditions, and actionable maintenance logic can be represented in the DiagnosticIQ framework. The case highlights cross-asset, multi-sensor reasoning a primary challenge captured by the DiagnosisIQPro dataset.

Figure 8 illustrates the schematic layout of a closed-loop water-cooled chiller system. The diagram captures the three key subsystems:

e 6: Industrial Asset Types Addressed in the Diagnostic Rules

972	
973	Table
974	Asset Type
975	AHU
976	AIIU
977	
978	Air Compresso
979	
980	Boiler
981	
982	Chiller
983	
984	Cooling Tower
985	
986	CDAC (Commo
987	CRAC (Compu
988	AC)
989	Fan
990	1 411
991	Heat Exchange
992	e
993	Plate & Frame
994	
995	Pump
996	
997	Terminal Unit
998	VAV (Vanial-1
999	VAV (Variable ume Unit)
1000	Condenser
1001	Condenser
1002	ERV (Energy
1003	Ventilator)
1004	Water Heater
1005	
1006	UPS (Unint
1007	Power Supply)
1008	Electrical Pane
1009	
1010	

Asset Type	Description
AHU	Conditions and circulates air as part of an HVAC system. Regulates airflow, temperature, and humidity in commercial buildings.
Air Compressor	Converts electrical or mechanical power into pressurized air for pneumatic systems and equipment.
Boiler	Heats water or other fluids for use in heating systems, industrial processes, or power generation.
Chiller	Removes heat from water using vapor-compression or absorption cycles; supplies chilled water to cooling systems.
Cooling Tower	Rejects heat from water-cooled systems by dissipating it into the atmosphere. Common in HVAC and process cooling.
CRAC (Computer Room AC)	Cools air in data centers to maintain safe temperature and humidity for IT equipment. Specialized HVAC component.
Fan	Drives air movement for ventilation, circulation, or cooling. Includes exhaust, supply, and return fans.
Heat Exchanger	Transfers heat between two fluids without mixing. Used for efficient thermal regulation in building systems.
Plate & Frame HX	A compact type of heat exchanger using stacked plates to transfer heat between fluid streams.
Pump	Moves liquids through mechanical force. Used in chilled water, hot water, and process fluid systems.
Terminal Unit	Regulates temperature and air delivery in individual building zones. Includes fan coil units and unit ventilators.
VAV (Variable Air Volume Unit)	Controls airflow to a zone by varying damper position, often part of demand-driven ventilation.
Condenser	Rejects heat from refrigerant cycles in chillers or heat pumps. Includes air- or water-cooled variants.
ERV (Energy Recovery Ventilator)	Transfers heat and moisture between exhaust and incoming fresh air streams to improve HVAC efficiency.
Water Heater	Provides domestic or process hot water, separate from large-scale boiler systems. May be gas or electric.
UPS (Uninterruptible Power Supply)	Supplies provides temporary backup power during grid interruptions to protect critical equipment.
Electrical Panel	Distributes power to facility circuits and equipment; may be monitored for load balancing or safety.

- Refrigerant Loop: Evaporator, compressor, condenser, and expansion valve arranged in a standard vapor-compression cycle.
- Chilled Water Loop: Chilled water pump circulates through the building's cooling coils and returns to the evaporator.
- Condenser Water Loop: Removes heat from the condenser and rejects it to the atmosphere via a cooling tower.

The Building Automation System (BAS) coordinates the entire process by issuing control signals to the compressor, pumps, and other critical components.

1020 1021 1022

1011 1012

1013 1014

1015

1016

1017

1018 1019

#### I.1 SYSTEM CONTEXT AND MOTIVATION

1023 1024

1025

Water-cooled chillers are used in high-performance HVAC systems where condenser heat is rejected via cooling towers. These systems rely on evaporators, compressors, condenser loops, and building automation systems (BAS) to coordinate thermal transfer. Rule-based diagnostics are essential for

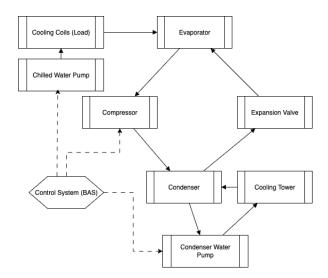


Figure 8: Schematic of a closed-loop water-cooled chiller system

early detection of inefficiencies or faults, particularly in critical environments like data centers and hospitals. As shown in Table 7, each component is associated with a distinct set of sensors and diagnostic KPIs.

Table 7: Components, Sensors, and KPIs in a Closed-Loop Water-Cooled Chiller System with Cooling Tower

Common Sen-	Associated KPIs	Purpose / Insight
sors / Meters		
Inlet Temp, Out-	$\Delta T_{\text{Evap}} = T_{\text{in}} - T_{\text{out}}$	Measures heat absorbed from
Flow	$\Delta T$	chilled water; identifies under- performance or fouling.
Power (kW),		Assesses mechanical load, effi-
Amps, Suc-		ciency, early signs of failure or
tion/Discharge	Compression Ratio =	degradation.
Pressure, Vibra-	$P_{ m dis}/P_{ m suc}$	
tion		
Inlet Temp, Out-		Evaluates heat rejection quality;
let Temp, Water		detects fouling, flow loss, scal-
Flow		ing.
	Superheat / Subcooling	Indicates refrigerant charge
Pressure	Temperatures	level, valve responsiveness, control precision.
Flow Pate	Pump Efficiency - Flow	Verifies circulation; detects
		pump wear, airlocks, or block-
		age.
	$\Delta T$ Stability	age.
	$\Delta T_{T_1,\ldots,T_n} = T_{T_1,\ldots,T_n} -$	Validates heat rejection; detects
		overcooling, bypass issues, fan
* / /	Approach Temp =	control faults.
		Control faults.
	2 supply 2 amolent	
	Setpoint Attainment =	Supports alarm generation, diag-
Status, Fault		nostics of control tuning, over-
Logs, Mode		shoot, inefficiency.
Indicators	Logs	, i
	Inlet Temp, Outlet Temp, Water Flow Power (kW), Amps, Suction/Discharge Pressure, Vibration Inlet Temp, Outlet Temp, Water Flow Pre/Post Temp, Pressure  Flow Rate, $\Delta P$ (Suction-Discharge), Power, Status Return Temp (from condenser), Supply Temp (to condenser), Fan Speed, OAT, Water Level Setpoints, Run Status, Fault Logs, Mode	Inlet Temp, Outlet Temp, Water Flow  Power (kW), Amps, Suction/Discharge Pressure, Vibration  Inlet Temp, Outlet Temp, Outlet Temp, Outlet Temp, Outlet Temp, Outlet Temp, Water Flow  Inlet Temp, Outlet Temp, Outlet Temp, Water Flow  Pre/Post Temp, Pressure  Flow  Rate, AP (Suction-Discharge), Power, Status  Return Temp (from condenser), Supply Temp (to condenser), Fan Speed, OAT, Water Level  Setpoints, Run Status, Fault Logs, Mode  Inlet Temp, Outlet $\Delta T_{\text{Cond}} = T_{\text{in}} - T_{\text{out}}$ $\Delta T_{\text{Cond}} = T_{\text{in}} - T_{\text{out}}$ Approach Temp = $T_{\text{refrigerant}} - T_{\text{cond-out}}$ Superheat / Subcooling Temperatures $\Delta P \text{ Stability}$ $\Delta T_{\text{Tower}} = T_{\text{return}} - T_{\text{supply}}$ Approach Temp = $T_{\text{supply}} - T_{\text{ambient}}$

#### I.2 COMPONENT-LEVEL DIAGNOSTIC RULES

We present a curated set of expert-authored diagnostic rules, each associated with a chiller subsystem and defined by time-persistent Boolean logic. Table 8 summarizes selected rules aligned with the DiagnosticIQ schema.

Table 8: Sample Diagnostic Rules for Closed-Loop Water-Cooled Chiller System

Component	Rule Name	Associated Sensors /	Condition Logic Summary
		KPIs	
Control Sys-	Cooling Temp	Supply Temp, Setpoint	Chiller Running \( \) (Supply
tem	Setpoint Not Met	Temp (KPI: $\Delta T_{\text{Setpoint}}$ )	Temp - Setpoint Temp $> 5^{\circ}F$ )
		-	for 2 hrs
Control Sys-	Low Supply Tem-	Supply Temp, Setpoint	Chiller Running \( \) (Setpoint -
tem	perature	Temp (KPI: $\Delta T_{\text{Setpoint}}$ )	Supply Temp $> 3^{\circ}F$ ) for 2 hrs
Evaporator	Cooling $\Delta T$ Low	Return Temp, Sup-	Chiller Running $\land$ P&F Off $\land$
		ply Temp, OAT (KPI:	$\Delta T < 7^{\circ} \text{F} \land \text{OAT} > 37^{\circ} \text{F} \text{ for}$
		$\Delta T_{ m Evap})$	4 hrs
Compressor	Efficiency Ex-	Power Input, Cooling	Chiller Running ∧ Efficiency >
	ceeds Threshold	Load (KPI: Chiller Effi-	design parameter for 2 hrs
		ciency)	
Condenser	Flow Detected	Condenser Water Flow,	Chiller Off $\land$ Flow $> 50$ GPM
	While Off	Run Status	for 2 hrs
Compressor	Load Low	Load %, Amps, Full	Chiller Running ∧ (Load % <
		Load Amps	$30\% \lor Amps / FLA < 30\%$ ) for
			2 hrs
Evaporator	Evaporator Ap-	Supply Temp, Refriger-	Chiller Running ∧ (Supply - Re-
	proach High	ant Temp (KPI: Evap Ap-	frigerant Temp $> 4$ °F) for 3 hrs
		proach)	
Condenser	Condenser Ap-	Liquid Temp, Return	Chiller Running ∧ (Condensate
	proach High	Temp (KPI: Cond Ap-	- Return Temp $> 4^{\circ}F$ ) for 3 hrs
		proach)	
Pump / Con-	$\Delta P$ Not at Set-	$\Delta P, \Delta P_{Setpoint}$	$ABS(\Delta P - \Delta P_{Setpoint}) > 4 PSI$
trol	point		for 3 hrs
Control Sys-	Excessive Power	Power, Run Status	Chiller Off $\land$ Power $> 5$ kW for
tem	While Off		3 hrs
Control Sys-	Chiller Cycling	Run Status Log	Status changed $\geq 4$ times in 8
tem			hrs

#### I.3 INTEGRATED DIAGNOSTIC SCENARIO

This case demonstrates cross-asset rule activation. The observed conditions span the chiller and cooling tower, requiring compositional reasoning. In the following, The meaning of OAT is Outside Air Temperature.

#### **Observed Conditions (QC):**

· Chiller Running

• Supply Temp - Setpoint Temp > 5°F

• Evaporator  $\Delta T < 7^{\circ}F$ 

• Condenser Water Return - Supply Temp < 5°F

• Cooling Tower Supply Temp < Setpoint - 4°F • OAT  $> 43^{\circ}F$ 

#### Activated Rules $(T_{R_i})$ :

• Chiller - Setpoint Not Met

- 1134 1135
- Chiller Cooling  $\Delta T$  Low
- 1136
- Cooling Tower Water Too Cold

impairing chiller performance.

1140

1141 1142 1143

1144

1145 1146

1147 1148 1149

1150 1151 1152

1153 1154

1155

1157

1156

1158 1159 1160

1161 1162 1163

1164 1165

1166

1167 1168 1169

1170 1171

1172 1173

1174

1175 1176 1177

1178 1179

1180 1181 1182

1183

1184

1185 1186

1187

**Observation**  $(O_i)$ : Cooling tower is overcooling condenser water, reducing head pressure and

**Action** ( $A_i$ ): Tune tower fan control logic or enable bypass to raise condenser water temperature.

This case study demonstrates the type of multivariate, cross-component diagnostic reasoning supported by the Diagnosis IQPro dataset. By grounding rule activation in realistic sensor logic and translating conditions into actionable maintenance decisions, the example highlights the practical value of structured condition-action QA benchmarks for industrial asset management.

#### DIAGNOSTIC RULE LOGIC CATEGORIZATION WITH EXAMPLES

This section presents a structured classification of diagnostic rules based on their underlying Boolean logic. Diagnostic rules are widely used in automated fault detection and energy analytics systems to evaluate sensor data from building systems such as air handling units (AHUs), chillers, boilers, and compressors. Understanding the logical structure of these rules enhances interpretability, facilitates rule development, and supports systematic debugging. Tables of 9 and 10 outline common logic categories and provide real-world examples to illustrate each structure.

Table 9: Boolean Logic Categories Used in Diagnostic Rules

Logic Category	Explanation
Conjunctive (AND)	All listed conditions must be simultaneously satisfied.
	Example pattern: $c_1 \wedge c_2 \wedge \cdots \wedge c_n$
Disjunctive (OR)	Any one condition is sufficient to trigger the rule.
	Example pattern: $c_1 \vee c_2 \vee \cdots \vee c_n$
Mixed (AND-OR)	Structured combinations of conjunctive and disjunctive logic, typi-
	cally in disjunctive normal form (DNF).
	Example pattern: $(c_1 \wedge c_2) \vee (c_3 \wedge c_4)$
Negation-based	Includes explicitly negated conditions.
	Example pattern: $c_1 \wedge \neg c_2$ or $\neg (c_1 \vee c_2)$

Categorizing diagnostic rules by their logic structure enables clearer reasoning about their behavior, performance, and possible interactions. It also aids in identifying redundant or conflicting rules in complex control systems. As building analytics platforms scale, such structured representations provide a foundation for more advanced techniques like rule validation, automated rule generation, and explainable diagnostics.

#### K DATASET GENERATION METHODOLOGY

**Algorithm Description:** The QA Generation Pipeline (See Algorithm 1) takes as input a set of expert-defined rules  $\{\mathcal{R}^i\}$ , asset descriptions (Desc), and a parameter max\_n\_choices controlling the maximum number of answer options per question. For each rule, the pipeline extracts atomic conditions from the condition tree  $TR^i$ , retrieves the corresponding asset description, and generates question-answer pairs by selecting and eliminating candidate observations using similarity metrics and heuristics. It then combines each extracted condition with all relevant question-option-answer tuples to build the final dataset  $DS_{\mathcal{O}}$ , facilitating systematic benchmarking of maintenance action recommendations.

#### K.1 RULE TO RULE SIMILARITY MAP

We utilize Rule to Rule Similarity **RRSim** mapping during the creation of dataset. The similarity is calculated by initially embedding the text components (asset\_type, conditions) of each rule using a all-mpnet-base-v2 embedding model. Then the embedding to embedding similarity is calculated according to cosine similarity.

1201 1202 1203

1205 1206

1226 1227

1228 1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

Table 10: Examples of Diagnostic Rules Categorized by Logical Structure.

Asset Name/Rule	Logic Category	Logical Expression Summary
Air Compressor - Pressure	Disjunctive (OR)	(ABS[Pressure – Setpoint] > 10 PSI OR
Setpoint Attainment	-	Pressure > 130 PSI)
AHU - Simultaneous Heat-	Mixed (AND-	AHU Running AND (Cooling Valve $\geq 5\%$
ing and Cooling	OR)	OR Preheat Valve $\geq$ 5%) AND (Drain Flags
		= 0)
AHU - Heating Valve Open	Mixed (AND-	AHU Running AND (OAT – SAT > 5°F OR
when Warm Outside	OR)	SAT Not Reporting) AND (Heating Valve >
		10% OR Preheat Valve > 10%)
Boiler - Excess O <sub>2</sub> in Stack	Disjunctive (OR)	(Gas Flow $>$ 5 AND Flue $O_2\%$ $>$ threshold
		OR Flue $O_2\%$ > threshold IF Fuel Flow Not
		Reporting)
CRAC - Limited Cooling	Conjunctive	CRAC Running AND (Return Temp ≤ Sup-
Warning	(AND)	ply Temp + 3°F)
Chiller - Cooling Substance	Conjunctive	Chiller Running AND (Supply Temp – Set-
Temperature Setpoint Attain-	(AND)	point $> 5$ °F)
ment		

#### Algorithm 1 QA Generation Pipeline

```
1207
              Input: \{\mathcal{R}^1, \dots, \mathcal{R}^{N_{\mathcal{R}}}\}, Desc, max\_n\_choices
1208
1209
              Output: DS_{\mathcal{O}}
1210
                1: Initialize DS_{\mathcal{Q}} \leftarrow []
1211
                2: for each \mathcal{R}^i \in \{\mathcal{R}^1, \dots, \mathcal{R}^{N_{\mathcal{R}}}\} do
                            \{QC_i^i\}_{i=1}^{N_{cond}} \leftarrow extracted\_conditions(\mathcal{TR}^i)
1212
1213
                            AD^i \leftarrow get\_asset\_desc(\mathcal{R}^i, Desc)
                4:
                             \{(QP_j^i, OPT_j^i, A_j^i)\}_{j=1}^{N_{sel}} \leftarrow extracted\_obs\_sel(\mathcal{R}^i, \alpha, RRSim, UO) 
 \{(QP_j^i, OPT_j^i, A_j^i)\}_{j=1}^{N_{eli}} \leftarrow extracted\_obs\_eli(\mathcal{R}^i, \beta, RRSim, UO) 
1214
                5:
1215
1216
                            all\_opts \leftarrow \{(QP_{j}^{i}, OPT_{j}^{i}, A_{j}^{i})\}_{j=1}^{N_{sel}} \cup \{(QP_{j}^{i}, OPT_{j}^{i}, A_{j}^{i})\}_{j=1}^{N_{eli}}
                7:
1217
                            for each QC_{j1}^i \in \{QC_j^i\}_{j=1}^{N_{cond}} do
                8:
1218
                                  for each (QP^i_{j2}, OPT^i_{j2}, A^i_{j2}) \in all\_opts do
                9:
1219
                                         Q^i \leftarrow (AD^i, QC^i_{j1}, QP^i_{j2}, OPT^i_{j2}, A^i_{j2})
1220
               10:
1221
                                         Append Q^i to DS_Q
               11:
1222
                                   end for
               12:
1223
               13:
                            end for
               14: end for
1224
               15: return DS_{\mathcal{O}}
1225
```

#### K.2 DIAGNOSISIQPRO

The **DiagnosisIQPro** dataset extends *DiagnosisIQ* to evaluate model performance under more challenging conditions with larger option sets. To enable direct comparison, for each question  $Q^i$  in Diagnosis Q, we retain the asset description  $AD^i$ , observed conditions  $QC^i$ , question prompt  $QP^i$ , and ground-truth answer  $A^i$ , while expanding the set of answer options  $OPT^i$  by adding additional plausible but incorrect choices.

For selection-type questions, we increase the number of incorrect options by resampling from observations of rules that are semantically similar yet distinct, ensuring that distractors remain relevant but incorrect. For elimination-type questions, we similarly augment the sets of correct and incorrect options to increase task complexity, leveraging domain-informed similarity measures (e.g., **RRSim**) to maintain logical coherence.

This augmentation more closely mimics real-world industrial scenarios, where practitioners must consider numerous potential failure causes, thereby testing the robustness and discriminative capabilities of language models in high-option environments.

#### K.3 DIAGNOSTICIQPERT

1242

1243 1244

1245

1246

1247

1248

1250 1251 1252

1279

1280

1282

1284

1285

1286

1287

We create a perturbation dataset DiagnosisIQPert to analyze the sensitivity of model responses to minor variations in the questions. This dataset is derived by manipulating DiagnosisIQ questions through several transformations: randomly shuffling the order of conditions and options, adding parentheses around option labels (e.g.,  $A \to (A)$ ), changing option labels (e.g.,  $A, B, C \to P, Q, R$ ), and substituting one question prompt  $QP^i$  with another. These perturbations help assess the robustness and consistency of language models when faced with slight changes in input formatting or phrasing.

## Figure 9: Example Question

```
1253
         Please select the correct option(s) from the following options given the question:
1254
         Ouestion:
1255
         ## Asset Description:
         AHU: Air Handling Unit: A device used to condition and
1256
         circulate air as part of a heating, ventilating, and air-conditioning (HVAC) system.
1257
1258
         ## Conditions:
1259
         - AHU Running
         - Outside Air Damper % < 15% AND Outside Air Damper Minimum % Not Reporting
         - Economizer Mode AND Supply Relative Humidity % Not Reporting
         - OAT < Setpoint Temperature
         - Outside Air Damper %
1263
         - OAT > 37 °F
         - Outside Air Damper % Does NOT = Daily Average
1265
         - SubType NOT OAU, RAS, RAU
1266
         ## How long the conditions were met:
1267
         Met for 2 Hours
1268
1269
         Looking at the current state of the asset, what is the MOST likely cause among the
1270
         options?
1271
         Options:
1272
         (P) Control system sent the wrong command
1273
         (0) Belts are loose or broken
1274
         (S) Broken Belt
         (R) Vanes at wrong angle
1275
         Your output must strictly follow this format:
1276
         {"answer": <the list of selected options, e.g., ["(P)", "(R)"]>}
```

#### K.4 DIAGNOSTICIQVERBOSE

To assess the symbolic understanding in the context of maintainance action recommendation we create a variant of DiagnosisIQ questions named DiagnosticIQVerbose. To create this we initially embedded the conditions of each question using a all-mpnet-base-v2 embedding model. Then the embedding are cluster according to cosine similarity to get 10 representative questions representing the cluster groups. We manually convert the conditions into natural language of these questions and then use these as in-context examples and prompt a mistral-large to generate the natural language representation for the rest of the questions.

```
1288
1289
       Your task is to read the asset description (## Asset Description:) and conditions (##
1290
            Conditions:) applied on the asset and
       write the conditions (## Conditions:) in natural language several examples are provided,
1291
             complete the last sample.
1292
1293
       ## Conditions:
1294
       AHU Running
1295
       OAT < 80F
      8 Cooling Valve % > 97%
```

1343

1344

1346 1347 1348

```
1296
      9 ABS(Supply Air Temperature Setpoint - Supply Air Temperature) > 3IF Setpoint Reporting
1297
     10
1298 11 ## Conditions in Natural Language:
1299 12 The asset is running while the outside temperatue is less than 80Fahrenheit and the
            units cooling valve is nearly fully open (Cooling Valve is open more than 97%)
1300
     13 and further tha absolute value of the difference between set threshold of air
1301
            temperature and supply air temperature is greater than 3
1302
     14 Fahrenheit
1303
     15
1304
    16
1305
    17
     18 ## Asset Description:
1306
     19 AHU: Air Handling Unit: A device used to condition and circulate air as part of a
1307
            heating, ventilating, and air-conditioning (HVAC) system.
1308 20
       ## Conditions:
1309 21
1310 22
       AHU Running
       OAT > 35F
     23
1311
       Preheat Valve % > 97%
     24
1312
1313 26
       ## Conditions in Natural Language:
1314
```

Listing 2: Prompt used to convert symbolic representation natural language

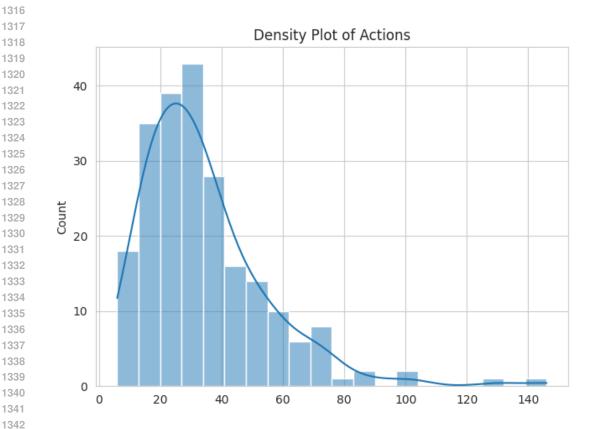


Figure 10: The word count distribution of unique actions in the expert curated dataset

#### K.5 QUESTION EXAMPLES

1352 1353

1350

1351

Figure 11: MAReE Prompt with 1 example question.

```
1354
         ## Asset Description:
1355
         Cooling Tower
1356
         ## Conditions:
1357
         - Cooling Tower Running
1358
         - 55 degF < Outside Air Temp < 80 degF
1359
         - Supply Temp Setpoint = Previous Hour Supply Air Temp Setpoint
         - Supply Temp Setpoint = Previous Daily Average Supply Air Temp Setpoint
1360
         ## How long the conditions were met:
1361
         Met for 3 Hours Checking Previous 3 Days Daily Average
         Analyse the given conditions of the presented asset and rank
1363
         the options that MOST likely gives the reason
         for the conditions?
1364
         A. Outside air temperature sensor failure
         B. Fans are off
1366
         C. Check fans and condenser water pumps
         D. VFD operation
1368
         E. Too many cnodenser pumps running
1369
         F. Logic issues for the cooling tower
         G. Cooling tower reset in manual
1370
         H. Fan is overridden
1371
         I. Too few cooling towers running
1372
         J. Static pressure sensors need calibration, repair or replacement
1373
         ## Use following Questions and answers as help for the ranking.
1374
         ### Example 1
         ### Asset Description:
1375
         Cooling Tower: A heat rejection device that cools water or other fluids
1376
         by transferring heat to the atmosphere. It is commonly used in HVAC systems,
1377
         power plants, and industrial processes.
1378
         ### Conditions:
         - Cooling Tower { Condenser Water is too cold
         - Cooling Tower Running
1380
         - OAT > 43 °F
1381
         - Condenser Water Supply Temperature to Chiller < 55 °F
         IF Condenser Water Temperature Setpoint NOT Reporting
         ### How long the conditions were met:
         Met for 2 Hours
         Review the listed conditions and identify which option MOST accurately accounts for
1385
1386
         A. Fan blades at incorrect pitch
1387
         B. Load is too low or fluctuates
1388
         C. Fan is overridden
         D. If unit resets based on VAV damper position exempt from this rule.
1389
         Answer: C. Fan is overridden
1390
         Your output must strictly follow this format:
1391
         {"option": <list of the option tag e.g. ['A', 'B', 'C', 'D', 'E']>,
1392
          score":<list of scoring value inline with rank ranging from 1,-1 eg: [1.0, 0.9, 0.8,
1393
         0.7, 0.6] >
1394
         "rank":<list of the rank eg: [10, 9, 8, 7, 6]>}
1395
         Your output in a single line:
1396
```

1398 1399

#### MAINTENANCE ACTION RECOMMENDATION ENGINE

1400 1401 1402

1403

The Maintenance Action Recommendation Engine takes a semi-defined rule as input in the form (asset\_type,  $TR^i$ ,  $t^i$ ), and returns a set of recommended maintenance actions  $O^i$  corresponding to a given rule  $\mathcal{R}^i$ .

1414 1415

1416

1417

1418 1419

1420

1421

1422

1423

1424 1425

1426

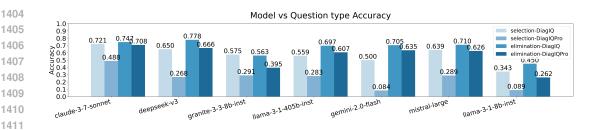


Figure 12: Model vs Question type Accuracy

We begin by selecting all possible actions applicable to a specific asset, resulting in a large action space AS. These actions are sourced from the collection of unique operations defined in our Expert-Curated Rule Documents (see Sec. 3.1).

To manage the scale of this action space, we adopt a divide-and-conquer strategy by chunking the actions into manageable segments, as described in Algorithm 2.

we utilize the DiagnosticIQ as a database of QA to inject examples giving the question and the answer rather

An example of a prompt generated by the get\_llm\_prompt function is shown in Fig. 11.

#### Algorithm 2 Maintenance Action Recommendation Engine

```
1427
           Input: DS_{\mathcal{O}}, asset_type, \mathcal{R}^i, N_{obs}, top_k
           Output: \mathcal{O}
1428
            1: function SELECTACTIONS(Q_{list}, R, obs, k)
1429
            2:
1430
                     \mathcal{Q}_{sel} \leftarrow \texttt{GET\_SIMILAR\_QUESTIONS}(\mathcal{Q}_{list})
            3:
                     prompt \leftarrow \text{GET\_LLM\_PROMPT}(\mathcal{R}, obs, \mathcal{Q}_{sel})
1431
            4:
                     order \leftarrow LLM\_ANSWER(prompt)
1432
            5:
                     Re-order obs using order
1433
            6:
                     sel\_obs \leftarrow \texttt{SELECT\_TOP\_OBSERVATIONS}(obs, k)
1434
            7:
                     return sel\_obs
1435
            8: end function
1436
            9: function DYNAMICACTIONRANKING(Q_{list}, Actions, asset_type, \mathcal{R}^i, N_{obs}, top_k)
1437
           10:
                     sel\_obs_{all} \leftarrow [\ ]
1438
           11:
                     for i \leftarrow 0 to |Actions| - 1 step N_{obs} do
1439
           12:
                          list\_obs \leftarrow Actions[i:i+N_{obs}]
                          sel\_obs \leftarrow SelectActions(Q_{list}, \mathcal{R}^i, list\_obs, top_k)
1440
           13:
           14:
                          Append sel\_obs to sel\_obs_{all}
1441
           15:
                     end for
1442
           16:
                     return DYNAMICACTIONRANKING(Q_{list}, sel\_obs_{all}, asset\_type, \mathcal{R}^i, N_{obs}, top_k)
1443
           17: end function
1444
               AS \leftarrow \text{GET\_ASSET\_ACTION\_SPACE}(asset\_type)
1445
           19: \mathcal{O} \leftarrow \text{DYNAMICACTIONRANKING}(DS_{\mathcal{O}}, AS, asset\_type, \mathcal{R}^i, N_{obs}, top_k)
1446
           20: return \mathcal{O}
1447
```