

EngiBench: A Benchmark for Evaluating Large Language Models on Engineering Problem Solving

Anonymous ACL submission

Abstract

Large language models (LLMs) have shown strong performance on mathematical reasoning under well-defined conditions. However, real-world engineering problems involve uncertainty, context, and open-ended settings that extend beyond symbolic computation. Existing benchmarks largely focus on well-defined or abstract reasoning and therefore fail to capture these complexities. We introduce EngiBench, a hierarchical benchmark designed to evaluate LLMs on solving engineering problems. It spans three levels of increasing difficulty (foundational knowledge retrieval, contextual reasoning, and open-ended modeling) and covers diverse engineering subfields. To facilitate a deeper understanding of model performance, we systematically rewrite each problem into three controlled variants (perturbed, knowledge-enhanced, and math abstraction), enabling us to separately evaluate the model’s robustness, domain-specific knowledge, and mathematical reasoning abilities. Experimental results show clear performance stratification across difficulty levels: model accuracy declines with task complexity, degrades under minor perturbations, and remains substantially below human performance on high-level engineering tasks. These findings reveal that current LLMs still lack the high-level reasoning needed for real-world engineering, highlighting the need for future models with deeper and more reliable problem-solving capabilities. Our source code and data are available at <https://anonymous.4open.science/r/EngiBench-2C7A>.

1 Introduction

Large language models (LLMs) have demonstrated promising capabilities in a range of mathematical reasoning tasks, from foundational skills such as basic computation and structured problem-solving (Cobbe et al., 2021), multi-step reasoning (Shao

et al., 2024; Wei et al., 2022), to more complex applications like mathematical modeling (Guo et al., 2025) and the generation or verification of mathematical proofs (Yang et al., 2023; Lin et al., 2025; Ren et al., 2025). However, just using mathematical reasoning is not enough for real-world applications. In practice, many applications arise not in abstract mathematical settings but in engineering contexts, where problems are grounded in physical systems and must handle uncertainty and real-world constraints. These characteristics require not only mathematical computation, but also need broader capabilities to understand engineering contexts and solve complex engineering problems.

Engineering problems differ fundamentally from mathematical problems (Hendrycks et al., 2021). Rather than seeking single closed-form answers, engineering requires finding feasible solutions that balance objectives under real-world constraints (Dym et al., 2005). For example, designing a drone system (Table 1) involves identifying operational requirements and managing trade-offs among range, payload, and energy limits. As shown in Figure 1, solving such problems requires more than recalling formulas or executing isolated calculations; it involves a sequence of interconnected cognitive steps, from understanding context and selecting appropriate assumptions to navigating trade-offs and addressing uncertainties. We refer to this broader set of competencies as the engineering problem-solving capability, consisting of four dimensions: *information extraction*, *domain-specific reasoning*, *multi-objective decision-making*, and *uncertainty handling*.

Despite the broader requirements of real-world engineering tasks, most existing benchmarks focus narrowly on well-defined mathematical problems. Benchmarks such as GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), and Omni-MATH (Gao et al., 2025) primarily assess symbolic reasoning, calculation, and formal problem-solving

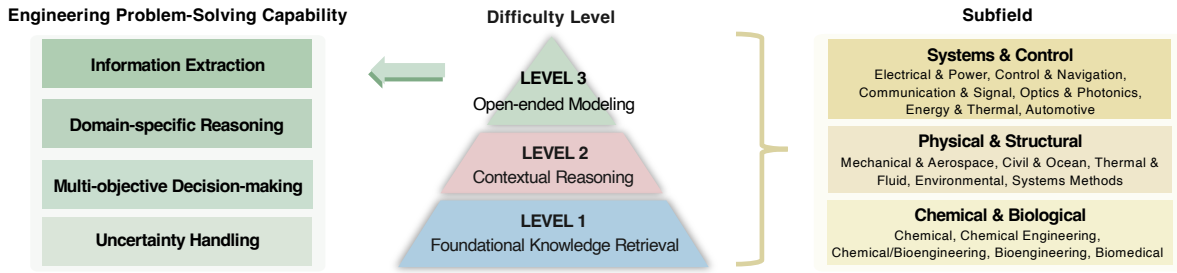


Figure 1: EngiBench taxonomy across capability dimensions, three difficulty levels, and three engineering subfields. Level 3 focuses on open-ended engineering reasoning.

under clean and well-defined conditions. Although some include basic engineering questions, they fail to capture the deeper reasoning required for real-world problem solving (Hendrycks et al., 2021; Wang et al., 2024c; Albalak et al., 2025; Du et al., 2025). A further limitation is that many benchmarks rely on public datasets without systematic rewriting, increasing the risk of pretraining overlap and inflated scores (Deng et al., 2024; Huang et al., 2025; Sainz et al., 2023). For example, GSM1k re-creates GSM8k-style questions to reduce overlap and observes performance drops of up to 8% (Zhang et al., 2024). Without such safeguards, evaluations may reflect memorization rather than genuine reasoning, providing limited insight into an LLM’s ability to address realistic engineering tasks.

In this work, we introduce EngiBench, an evaluation framework designed not only to assess LLMs’ engineering problem-solving capabilities but also to diagnose where and why these capabilities fail. The benchmark spans multiple engineering subfields and structures tasks into three difficulty levels that reflect the progression from foundational knowledge retrieval to contextual reasoning and open-ended modeling. To support fine-grained diagnosis, each problem is provided in three controlled variants that separate robustness, domain knowledge, and mathematical reasoning. Evaluation centers on four capability dimensions essential to engineering problem solving: *information extraction*, *domain-specific reasoning*, *multi-objective decision-making*, and *uncertainty handling*. For open-ended tasks, we further adopt rubric-based evaluation using expert-designed scoring criteria to ensure consistent and reliable assessment. Together, these components create a diverse, high-quality, and contamination-aware benchmark for evaluating LLMs’ engineering problem-solving capabilities.

Experiment results show clear stratification across difficulty levels, with higher-level tasks high-

lighting distinct capability gaps. The perturbed variant leads to performance drops, even in strong models, revealing that prior evaluations may overestimate true generalization. Most importantly, current LLMs perform poorly on Level 3 tasks involving open-ended, high-level engineering reasoning and remain far below human experts. These findings suggest that today’s LLMs are still far from reliably addressing real-world engineering problems, leaving substantial room for future improvement.

Our contributions can be summarized as follows: (1) We are among the first to systematically evaluate LLMs on real-world engineering problems; (2) We design a hierarchical benchmark with three difficulty levels and multiple problem variants, enabling fine-grained analysis of model reasoning capabilities and limitations; (3) Unlike prior benchmarks, our benchmark systematically evaluate LLM performance on open-ended engineering tasks; (4) We evaluate a broad set of mainstream LLMs, providing insights that can aid future model development and enhance engineering capabilities.

2 Related Works

LLMs for Engineering Problems. LLMs integrate broad domain knowledge with strong multi-step reasoning, making them promising tools for complex problem solving. Engineering problems, however, require modeling real-world systems and reasoning under practical constraints. Despite the growing use of LLMs in engineering applications, their true engineering problem-solving capability remains unclear due to the limitations of existing benchmarks (Wang et al., 2024b; Ma et al., 2024; Tang et al., 2024; Cheng et al., 2025). General-purpose benchmarks, including MMLU (Hendrycks et al., 2021), MMLU-Pro (Wang et al., 2024c), BIG-Math (Albalak et al., 2025), and SuperGPQA (Du et al., 2025), contain only limited engineering content. Most questions emphasize factual recall in multiple-choice form, fail-

Table 1: Hierarchical difficulty from mathematics to real-world engineering. This illustrates three levels of increasing complexity. Examples show the progression from closed-form math problems to open-ended engineering scenarios.

Level	Definition	Example
Mathematics	Mathematical tasks are typically well-posed and self-contained , with complete information and clearly defined solution spaces.	A machine produces 45 parts per minute. If it operates continuously for 2 hours, how many parts will it produce in total? <i>⚡ This task requires only basic multiplication and does not involve any domain knowledge. It represents a typical closed-form numerical computation problem.</i>
Upgrading Condition: Incorporating domain-specific engineering knowledge		
Engineering Level 1: Foundational Knowledge Retrieval	Apply basic engineering concepts or formulas to structured problems via single-step computation.	A drone operates at a constant power of 200W for 30 minutes. Calculate the total energy consumption in joules. <i>⚡ This task requires applying the basic physical formula $E = P \times t$, with unit conversion from minutes to seconds. It tests the model's ability to retrieve and apply foundational engineering knowledge in a single-step calculation.</i>
Upgrading Condition: Multi-step reasoning and contextual integration		
Engineering Level 2: Contextual Reasoning	Perform multi-step reasoning under well-defined constraints by integrating conditions and domain knowledge.	A drone needs to fly 6 km. The first half is uphill, increasing power usage by 20%, while the second half is flat at 180W. The drone flies at 30 km/h and uses a battery rated at 8000mAh, 11.1V. Can the battery support the trip? <i>⚡ This task requires multi-step reasoning: estimate flight time, adjust power consumption, and compare with battery capacity.</i>
Upgrading Condition: Solving open-ended, under-specified problems		
Engineering Level 3: Open-ended Modeling	Solve open-ended , real-world problems through information extraction, trade-off reasoning, and uncertainty handling.	Design a drone system for urban delivery that balances multiple factors, including flight range, payload capacity, and cost control. Propose a feasible solution and justify your design decisions. <i>⚡ This is an open-ended problem with incomplete constraints and potentially conflicting objectives, requiring information extraction, trade-off analysis, and robustness under uncertainty.</i>
🔍 Information Extraction	Identify and extract relevant information from complex or redundant problem descriptions.	Identify critical variables —such as payload weight, wind speed, flight duration, and battery margin—from complex or verbose task descriptions.
📚 Domain-specific Reasoning	Apply specialized engineering principles and structured knowledge to guide logical inference and solution formulation.	Apply specialized engineering knowledge —such as flight mechanics and battery discharge principles—to formulate models and perform technical analysis.
🎯 Multi-objective Decision-making	Make justified trade-offs between competing in the absence of a single optimal solution.	Justify trade-offs among competing objectives like range, cost, safety, and operational efficiency when no single optimal solution exists.
🌀 Uncertainty Handling	Ensure solution robustness by reasoning under incomplete, variable, or ambiguous real-world conditions.	Account for unpredictable factors such as weather, task variation, and battery aging, and design robust strategies (e.g., adding 20% battery reserve) to ensure reliable performance.

ing to capture core engineering reasoning. Several domain-specific engineering benchmarks have been proposed, including EEE-Bench (Li et al., 2024), ElecBench (Zhou et al., 2024), FEABench (Mudur et al., 2025), TransportBench (Syed et al., 2024), and JEEBench (Arora et al., 2023). However, they largely focus on single disciplines and well-defined tasks, limiting their ability to evaluate open-ended and cross-disciplinary engineering reasoning. To address this gap, we introduce a multi-level engineering benchmark spanning multiple subfields and incorporating both closed-form and open-ended tasks, enabling a comprehensive evaluation of engineering capabilities.

LLM for Mathematical Problems. Mathematics has been widely used to evaluate LLMs because it requires structured logic, multi-step deduction, and rigorous formal reasoning. Existing benchmarks cover elementary problems (Cobbe et al., 2021; Hendrycks et al., 2021; Patel et al., 2021; Amini et al., 2019), advanced symbolic reasoning (Hendrycks et al., 2021; Albalak et al., 2025), theorem proving (Zheng et al., 2022; Liang et al., 2023; Gao et al., 2025; Lu et al., 2024), and multi-modal tasks (Wang et al., 2024a). However, these benchmarks do not capture engineering-specific reasoning, such as system modeling, reasoning under constraints, or domain assumptions. Our work builds on their evaluation and shifts the focus to real-world engineering tasks.

Evaluation Challenges. Evaluating LLMs on engineering problem solving is challenging due to the complexity and open-ended nature of engineering tasks. Existing evaluation approaches can be broadly grouped into reference-based, task-

oriented, preference-based, and rubric-based methods. Reference-based and task-oriented evaluations work well when ground truths or executable outputs are available, as in MathVista (Lu et al., 2024), CHAMP (Mao et al., 2024), EEE-Bench (Li et al., 2024), and FEABench (Mudur et al., 2025). However, many core engineering capabilities cannot be captured in such closed-form settings. Preference-based methods, such as MT-Bench-101 (Bai et al., 2024), use pairwise comparisons for open-ended tasks but are often influenced by model-specific generation patterns, limiting objectivity. Rubric-based evaluations score responses along multiple criteria, with general-purpose frameworks like Prometheus (Kim et al., 2024) focusing on abilities such as context retention and paraphrasing rather than engineering reasoning.

3 Methodology

3.1 Engineering Problem-Solving Capability

Engineering problems require context-aware solutions under real-world constraints, distinguishing them from mathematical problems that typically operate in well-defined, closed-form settings (Dym et al., 2005; Hendrycks et al., 2021). Beyond abstraction and logical rigor, engineering problem solving involves a sequence of interconnected cognitive steps, from interpreting problem context to making decisions under constraints and uncertainty. We refer to this as engineering problem-solving ability, which comprises four key dimensions: information extraction, domain-specific reasoning, multi-objective decision-making, and uncertainty handling. These dimensions align with established paradigms in engineering modeling, in-

cluding information filtering, constraint-based and multi-objective formulation, and robustness analysis (see Figure 1 and Table 1).

Information Extraction. The capability to identify and organize key variables, constraints, and objectives from complex or noisy descriptions. It reflects the model’s capacity to filter irrelevant information and transform unstructured text into structured representations for downstream reasoning.

Domain-specific Reasoning. The capability to apply engineering knowledge, including physical principles, empirical rules, and domain conventions, to interpret scenarios and choose appropriate solution strategies. It involves recognizing valid approximations, implicit assumptions, and methods used in engineering practice.

Multi-objective Decision-making. The capability to balance competing objectives such as cost, performance, and safety when no single optimal solution exists. This dimension assesses whether a model can justify trade-offs under constraints, a defining feature of engineering problem solving.

Uncertainty Handling. The capability to reason under incomplete, noisy, or dynamic information. It includes anticipating uncertainties, incorporating safety margins or fallback strategies, and generating solutions that remain robust despite ambiguity. This capability is essential for making reliable engineering decisions in real-world settings.

3.2 Problem Hierarchical Difficulty Design

Engineering problem solving involves multiple distinct capabilities, making it difficult to assess through any single task or a one-dimensional hierarchy. A clear taxonomy is therefore essential for identifying where models succeed or fail. To provide such structure, EngiBench organizes engineering tasks into three complementary levels: foundational knowledge retrieval, contextual reasoning, and open-ended modeling, each reflecting different cognitive demands. Rather than simply aggregating tasks, this framework aligns evaluation with the core capabilities, and its effectiveness is demonstrated in Section 4.2.1.

Level 1. Tasks are well-defined and self-contained, typically requiring only single-step application of fundamental engineering formulas. They emphasize factual recall, accurate computation, and minimal contextual reasoning. This level assesses whether a model has a stable engineering knowledge base and can reliably retrieve and apply it to straightforward problems.

Level 2. Tasks require multi-step reasoning under contextual constraints such as units, physical limits, and coupled variables. Although these problems are well-defined and have unique solutions, models must interpret structured descriptions and integrate domain knowledge across steps to generate correct answers. Compared with Level 1, simple recall is insufficient; models need to handle structured complexity to generate correct solutions.

Level 3. Tasks reflect open-ended engineering challenges with uncertainty, incomplete information, and conflicting objectives. They require the full engineering problem-solving capability. Unlike Level 1 and Level 2, problems do not have a single correct answer, and evaluation focuses on how well models demonstrate robust and adaptive reasoning under open-ended conditions.

3.3 Dataset Construction

Data Sources. We collect data from three primary sources: problems selected from existing public benchmarks, university educational materials, and modeling competitions. These problems reflect the intended hierarchy of difficulty described above and address the lack of open-ended engineering modeling problems with expert-defined evaluation criteria in existing datasets.

Construction Process. Levels 1 and 2 contain structured engineering problems with standard answers, drawn from general-domain benchmarks such as SuperGPQA (Du et al., 2025), MMLU (Hendrycks et al., 2021), MATH (Hendrycks et al., 2021), GSM8k (Cobbe et al., 2021), Orca-Math (Mitra et al., 2024), HARP (Yue et al., 2024), Omni-MATH (Gao et al., 2025), Big-Math (Albalak et al., 2025), and selected university resources. Although these datasets are broad in scope, all problems used in EngiBench were passed through an engineering relevance filtering procedure (Appendix D.1) to retain only questions that align with engineering knowledge. All selected problems were further standardized and validated.

Level 3 introduces the first systematic collection of open-ended engineering tasks, comprising 43 problems from major modeling competitions. Each problem includes official scoring rubrics and reference solutions provided by top-ranking competition winners. All task rewrites and scoring rubrics were finalized by domain experts, with LLMs providing auxiliary support during intermediate steps, ensuring clarity, rigor, and reliable assessment (see Appendix D.2 and Appendix F.2).

Problem Annotation and Quality Control.

Level 3 problems and their scoring rubrics were expert-reviewed by 20 PhD students and engineering professionals, with LLMs used only as auxiliary tools. From nearly 1,000 competition questions, we retained only those with official rubrics and performed extensive text-based reformulation of formulas, tables, and diagrams. The released scoring scripts implement these expert-defined rubrics for reproducible downstream evaluation. All annotation and quality control in this section are limited to problem and rubric construction (see Appendix F.2.1), while model response scoring is described in Section 4.1.

Coverage and Classification. EngiBench spans three subfields: Systems & Control (939 problems), Physical & Structural (354 problems), and Chemical & Biological (467 problems). This categorization reflects differences in problem focus, underlying domain knowledge, and the reasoning processes required to solve them.

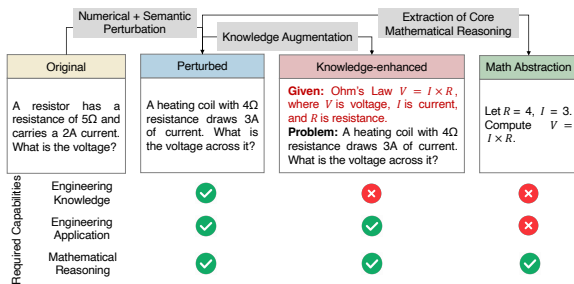


Figure 2: Each problem is rewritten into three controlled variants, each isolating a different reasoning capability.

3.4 Controlled Problem Variants

Evaluating LLMs on engineering tasks requires more than measuring overall accuracy. A correct answer may arise from data memorization rather than reasoning (Huang et al., 2025; Zhang et al., 2024; Mirzadeh et al., 2025; Srivastava et al., 2024; Gulati et al., 2024), while an incorrect answer may reflect missing domain knowledge, weak mathematical skills, or failures in interpreting engineering constraints. Without separating these factors, accuracy alone provides limited diagnostic value.

To enable deeper analysis, each problem is rewritten into three controlled variants derived from the original form (Figure 2). (1) The *perturbed variant* introduces numerical and semantic changes to assess robustness and reduce possible overlap with pretraining data. (2) The *knowledge-enhanced variant* adds essential domain information such as for-

mulas, constants, and key definitions so that errors caused by missing knowledge can be distinguished from reasoning failures. (3) The *math abstraction variant* removes contextual and domain-specific elements while preserving the underlying mathematical structure, allowing us to isolate mathematical reasoning and quantify how much engineering context affects performance.

These controlled variants provide a structured way to distinguish why a model succeeds or fails, giving a capability-oriented evaluation of engineering problem solving. For Level 1 and Level 2, all three variants are constructed systematically from the original problem. For Level 3, the open-ended nature and inherent complexity make knowledge-enhanced and math abstraction variants impractical, so only the perturbed variant is included.

4 Experiments

4.1 Experiment Setup

Evaluated LLMs. As the first batch, 16 LLMs were evaluated under the zero-shot setting, covering a representative range of model types. Specifically, we include: (1) closed-source models such as GPT-4.1, GPT-4.1 Mini, and GPT-4.1 Nano from OpenAI (Achiam et al., 2023); Claude 3.7 Sonnet and Claude 3.5 Sonnet from Anthropic (Anthropic, 2024b,a); and Gemini 2.5 Flash and Gemini 2.0 Flash from Google DeepMind (Team et al., 2023, 2024); (2) open-source models, including GLM-4-32B and GLM-4-9B from THUDM (GLM et al., 2024), Qwen2.5-72B and Qwen2.5-7B from Alibaba (Yang et al., 2024), Llama 4 Maverick (referred to as Llama 4) and Llama 3.3-70B (referred to as Llama 3.3) from Meta (Grattafiori et al., 2024), and DeepSeek-V3-671B (referred to as DeepSeek-V3) and DeepSeek-R1-Distill-Qwen-1.5B (referred to as DeepSeek-R1 7B) from DeepSeek (Liu et al., 2024; Guo et al., 2025), Mixtral-8x7B-Instruct-v0.1 (referred to as Mixtral 8x7B) from Mistral AI (Jiang et al., 2024). This selection spans diverse model sizes, training paradigms, and accessibility levels. We ensured consistent formatting and output parsing across all models.

Evaluation protocols. Level 1 and Level 2 consist of well-defined problems with unique solutions and are evaluated using binary scoring. Evaluation consistency is verified through multi-model cross-checking and random human spot checks. Further details are provided in Appendix F.1. Level 3 tasks are open-ended and are evaluated using a rubric-

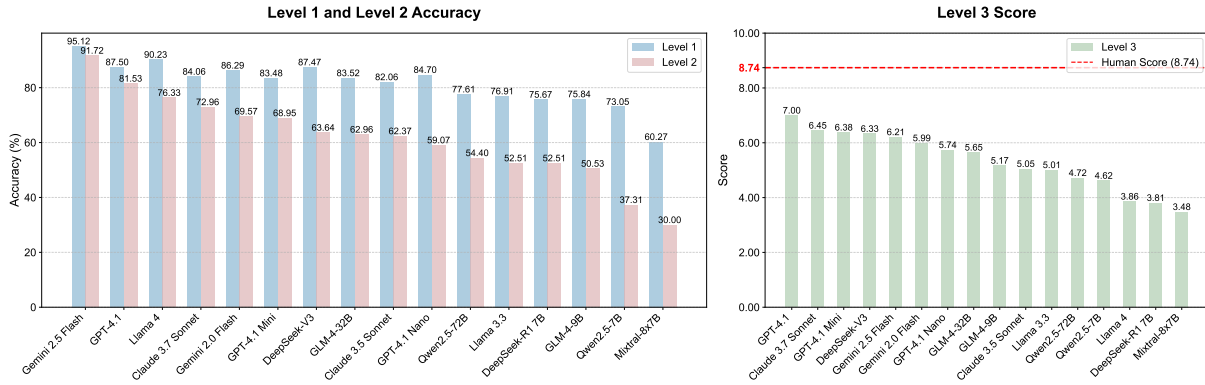


Figure 3: Overview of model performance across engineering reasoning tasks. The left subfigure shows model accuracy on Level 1 and Level 2 tasks, while the right subfigure presents scores on Level 3 open-ended tasks, with the human expert score indicated by the red line.

based framework derived from official criteria and refined by domain experts. Scoring is performed by LLMs following the same rubrics, and all Level 3 scores are subsequently reviewed and calibrated by human annotators following the same criteria. Further details are provided in Appendix F.2.2 and Appendix G.1.

Also, we introduce human scores for Level 3 tasks for comparison with LLMs’ performance. We obtain human scores from two sources: award-winning competition submissions (original version) and manual solutions by top-performing students for the perturbed variant. All human and LLM responses are evaluated using the same rubric to ensure consistency and fairness.

4.2 Results

4.2.1 Overall

Model stratification and design validation. Model performance exhibits a clear downward trend from Level 1 to Level 3, demonstrating the effectiveness of our hierarchical difficulty design. As shown in Figure 3, most models achieve high accuracy on Level 1, perform moderately on Level 2, and show a clear performance decline on Level 3. This trend indicates that our hierarchical framework successfully separates problems by cognitive difficulty, with each level reflecting distinct capability thresholds. The results validate that a multi-level design is necessary to capture the full range of engineering problem-solving capabilities.

Evaluating high-level engineering reasoning. Level 3 is designed to assess high-level engineering reasoning that goes beyond formulaic computation. Unlike Level 1 and Level 2, which focus on structured problem solving, Level 3 features open-ended and underspecified tasks that better reflect real-world engineering challenges. The sharp

performance drop at this level reveals the current limitations of LLMs in handling such complex scenarios. Besides, the gap between LLMs and human experts at Level 3 also reveals a key deficiency in high-level engineering capabilities. All evaluated models score well below the human expert, who achieves an average of 8.74, indicating that current LLMs are still far from reliably handling complex engineering problems. This underscores the need for further research to bridge this gap.

Smaller-scale LLMs struggle with complex tasks. While all LLMs show room for improvement on complex, open-ended engineering tasks, smaller-scale LLMs exhibit significantly greater limitations. As task complexity increases, performance disparities widen. At Level 1, most models still cluster within 70–90%. But at Level 2, leading models such as GPT-4.1 and Gemini 2.5 Flash achieve accuracies above 80%, whereas DeepSeek-R1 7B reaches only about 52% and other lightweight models often fall below 40%. This divergence is most evident at Level 3, where state-of-the-art models approach scores of 7.0, while lightweight models remain under 4.0. These results show that EngiBench is not saturated and continues to distinguish models across scales.

Robustness and contamination risk. Some LLMs may achieve high scores not through internal reasoning, but due to overlap with pretraining data. To reveal this, we use perturbed variant that apply minor contextual and numerical changes but keep the core structure unchanged. As shown in Figure 4, model performance remains relatively stable on Level 1 but drops sharply on Level 2. For example, on Level 2, accuracy decreases by 9.3% for GPT-4.1 Nano, 11.4% for Qwen2.5-7B, and 8.3% for Mixtral-8x7B. These declines suggest a stronger reliance on surface-level pattern match-

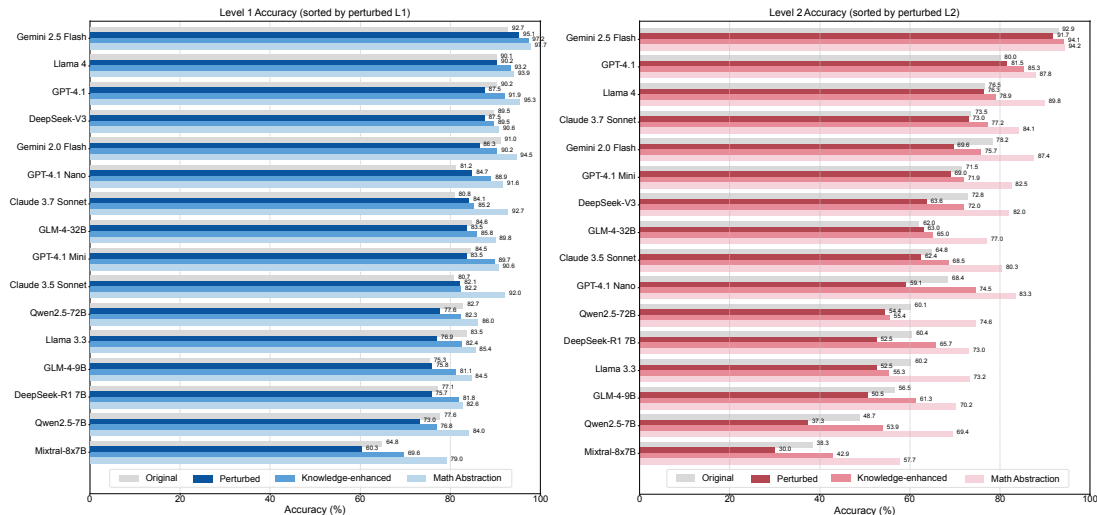


Figure 4: Accuracy of LLMs on Level 1 (left) and Level 2 (right) across the Original, Perturbed, Knowledge-enhanced, and Math Abstraction variants. Drops under the Perturbed variant reflect sensitivity to input changes, while gains on the latter two indicate that models benefit from added knowledge or simplified formulations.

ing, rather than robust reasoning, highlighting the role of perturbation-based evaluation in diagnosing overestimated capabilities.

4.2.2 Performance for Level 1 & Level 2 Tasks

Knowledge Enhancement Improves Accuracy.

Adding explicit domain knowledge consistently improves accuracy across all levels, especially for weaker models. As shown in Figure 4, models perform better on knowledge-enhanced variants than on perturbed inputs. These gains suggest two main failure sources: lacking essential domain knowledge or failing to apply it correctly during reasoning. Providing explicit knowledge therefore offers a clear diagnostic signal that helps distinguish knowledge deficits from reasoning errors, which is a key capability for engineering evaluation.

Math Abstraction Further Improves Performance.

LLMs perform even better when engineering problems are rewritten into purely mathematical form, removing contextual details. As shown in Figure 4, most models achieve their highest accuracy under this variant, especially smaller models that struggle with contextual interpretation. This pattern suggests that the main challenge in engineering tasks is not computation, but the earlier step of translating natural-language descriptions into well-structured mathematical formulations. This underscores the importance of evaluating the reasoning steps that precede formula application, as these upstream processes are not captured by traditional math benchmarks.

Smaller Models Are More Sensitive to Input Variants.

Smaller-scale LLMs exhibit much larger performance fluctuations across input ver-

sions, indicating limited generalization and unstable reasoning. As shown in Figure 4, in Level 2, Qwen2.5-7B drops by 11.4% under the perturbed variant, yet gains 16.6% with added domain knowledge and another 15.5% under math abstraction. In contrast, Gemini 2.5 Flash remains highly stable: its accuracy decreases by only 1.2% under the perturbed version and increases by 2.4% and 2.5% under the knowledge-enhanced and math abstraction variants, respectively. This comparison shows that smaller models are more sensitive to input formulation and tend to rely on surface patterns rather than consistent, context-aware reasoning.

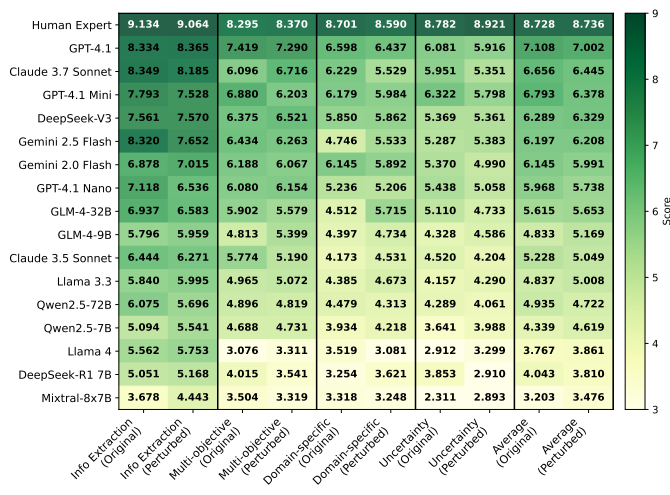
4.2.3 Performance for Level 3 Tasks

Dimension-wise and model-wise performance.

As shown in Figure 5a, human experts lead across all four dimensions with a balanced capability profile. In contrast, LLMs show uneven performance across the four dimensions. They handle redundant information extraction relatively well and perform moderately on multi-objective decision-making but struggle with domain-specific reasoning and uncertainty handling. This pattern indicates that their abilities are imbalanced, with clear deficiencies in key engineering-oriented skills. Results also demonstrate that model performance correlates with scale and accessibility. Larger, closed-source models like GPT-4.1 and Claude 3.7 Sonnet, consistently achieve average scores above 6. In contrast, smaller open-source models (e.g., Mixtral-8x7B) average below 4, with common omissions in aspects such as trade-off reasoning and uncertainty consideration.

Correlation analysis.

To quantify this trend, Fig-



(a) Level 3 Model Evaluation.

Information Extraction	Multi-objective Decision-making
Selection of Evaluation Indicators (6 pts) <ul style="list-style-type: none"> 6 pts: Covers efficiency, safety, robustness; clear formulas provided 4 pts: Includes reasonable indicators, but lacks full coverage or definitions 2 pts: Incomplete or loosely relevant indicators 0 pts: No valid indicators proposed 	Multi-Objective Optimization (6 pts) <ul style="list-style-type: none"> 6 pts: Formal multi-objective model (e.g., efficiency vs. safety vs. robustness) 4 pts: Mentions trade-offs but lacks full model 2 pts: Only single-objective considered 0 pts: No mention of optimization
Assumption Analysis (4 pts) <ul style="list-style-type: none"> 4 pts: Assumptions clearly stated and justified 2 pts: Lists assumptions, but lacks analysis 0 pts: No assumptions, or assumptions are irrelevant 	Computational Efficiency (4 pts) <ul style="list-style-type: none"> 4 pts: Efficient model; supports multiple scenario simulations 2 pts: Model works but inefficient 0 pts: No mention of runtime or efficiency
Uncertainty Handling	Domain-specific Reasoning
Modeling Traffic Variability (6 pts) <ul style="list-style-type: none"> 6 pts: Models peak/off-peak flows or stochastic variation 4 pts: Mentions variability, lacks modeling detail 2 pts: Weak or vague handling of uncertainty 0 pts: Ignores uncertainty 	Application of Traffic Flow Theory (5 pts) <ul style="list-style-type: none"> 5 pts: Correct use of flow-density-speed relationships or queuing theory 3 pts: Partial or incorrect theory use 0 pts: No use of traffic theory
Risk Evaluation & Mitigation (4 pts) <ul style="list-style-type: none"> 4 pts: Provides risk assessment and detailed response strategy 2 pts: Mentions risk, lacks concrete measures 0 pts: No discussion of risk 	Urban Planning & Traffic Management (5 pts) <ul style="list-style-type: none"> 5 pts: Proposes actionable, planning-based recommendations 3 pts: General suggestions not tied to planning 0 pts: No practical recommendations

(b) Scoring rubric example.

Figure 5: Level 3 Model Evaluation and Scoring Rubric. This figure summarizes Level 3 evaluation results and scoring standards. Subfigure (a) reports average model scores across four capabilities under both original and perturbed inputs. Subfigure (b) shows an example rubric outlining scoring criteria across capability dimensions.

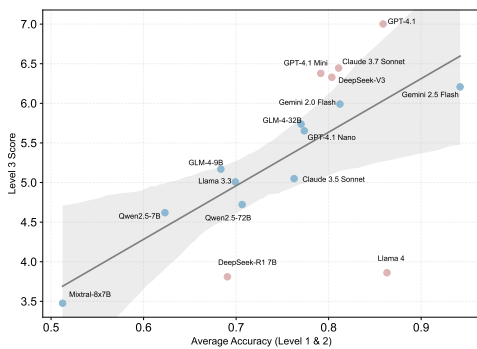


Figure 6: Correlation between structured tasks (Level 1&2) and open-ended tasks (Level 3).

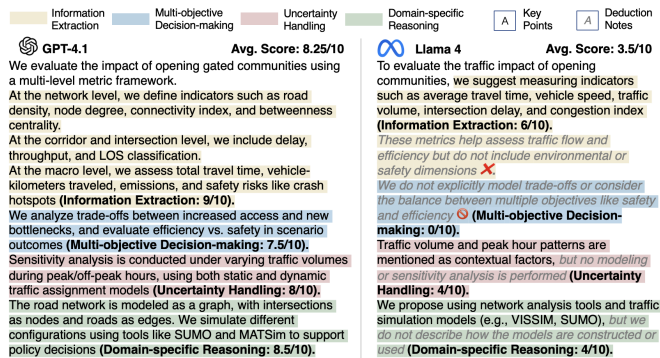


Figure 7: Case study showing why Llama 4 received low Level 3 scores.

Figure 6 illustrates the relationship between model performance on structured tasks (Levels 1 & 2) and open-ended tasks (Level 3). Overall, we observe a clear positive correlation: *models that achieve higher accuracy on structured tasks tend to also perform well on open-ended tasks*, suggesting a general consistency across task types. At the same time, some models deviate from this general trend. GPT-4.1, Claude 3.7 Sonnet, and DeepSeek-V3 show notably stronger performance on Level 3 than their results on Levels 1 and 2 would suggest, indicating more advanced reasoning and modeling abilities than what structured tasks alone reveal.

In contrast, models like Llama 4 perform pretty well on structured tasks but falter on open-ended ones, revealing weak high-level reasoning. Figure 7 illustrates this gap: Llama 4 scores 0 in multi-objective decision-making due to missing trade-off analysis, while GPT-4.1 provides a structured evaluation and scores 7.5. A similar shortfall also

appears in uncertainty handling. These examples show that Llama 4 can recall facts but struggles to apply them in complex, judgment-based scenarios.

5 Conclusion

We introduce **EngiBench**, a benchmark for evaluating LLMs on engineering problem solving across increasing levels of complexity. Our results show that while current models perform well on foundational knowledge retrieval, their performance declines significantly in multi-step contextual reasoning tasks, due to both domain knowledge gaps and limited mathematical reasoning. On open-ended modeling tasks, even the strongest models fall short of human-level performance, revealing persistent limitations in high-level reasoning, trade-off analysis, and uncertainty handling. These findings underscore the need for LLMs to move beyond pattern matching and toward deeper reasoning capabilities for real-world engineering applications.

6 Limitations

While EngiBench provides the first systematic evaluation of LLMs on real-world engineering problems, covering multi-level tasks, variant-based reasoning diagnostics, and open-ended modeling, several limitations remain that we plan to address in future work.

Multimodal Support. Many real-world engineering problems involve visual elements such as diagrams, schematics, or structured tables. The current version of EngiBench does not include multimodal tasks, as most existing LLMs still lack stable and consistent multimodal input capabilities. To avoid confounding engineering reasoning performance with visual processing variability and to ensure fair and comparable evaluation across models, we restrict all inputs to text-only formats.

Long-Context Support. Some engineering tasks involve long problem descriptions or extensive tabular data that exceed the input length limits of current LLMs. To avoid unfair model truncation effects and ensure uniform evaluation settings, such problems are not included in this version of the benchmark.

Human-in-the-loop Construction. Building the dataset involves substantial human effort, including problem collection, answer generation, and variant validation. This ensures data quality and alignment with engineering standards, but also reflects the significant manual effort behind the benchmark.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Alon Albalak, Duy Phung, Nathan Lile, Rafael Rafailov, Kanishk Gandhi, Louis Castricato, Anikait Singh, Chase Blagden, Violet Xiang, Dakota Mahan, and 1 others. 2025. Big-math: A large-scale, high-quality math dataset for reinforcement learning in language models. *arXiv preprint arXiv:2502.17387*.

Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*.

Anthropic. 2024a. **Claude-3 family: Opus, sonnet, haiku.** Available at: [https://](https://assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf)

assets.anthropic.com/m/61e7d27f8c8f5919/original/Claude-3-Model-Card.pdf. 657
658

Anthropic. 2024b. **Claude-3.5 sonnet.** Available at: <https://www.anthropic.com/news/claude-3-5-sonnet>. 659
660
661

Daman Arora, Himanshu Singh, and 1 others. 2023. Have llms advanced enough? a challenging problem solving benchmark for large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7527–7543. 662
663
664
665
666
667

Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and 1 others. 2024. Mt-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. *arXiv preprint arXiv:2402.14762*. 668
669
670
671
672
673

Yuheng Cheng, Huan Zhao, Xiyuan Zhou, Junhua Zhao, Yuji Cao, Chao Yang, and Xinlei Cai. 2025. A large language model for advanced power dispatch. *Scientific Reports*, 15(1):8925. 674
675
676
677

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*. 678
679
680
681
682
683

Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gestein, and Arman Cohan. 2024. **Investigating data contamination in modern benchmarks for large language models.** In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8706–8719, Mexico City, Mexico. Association for Computational Linguistics. 684
685
686
687
688
689
690
691
692

Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, and 1 others. 2025. Supergpqa: Scaling llm evaluation across 285 graduate disciplines. *arXiv preprint arXiv:2502.14739*. 693
694
695
696
697

Clive L Dym, Alice M Agogino, Ozgur Eris, Daniel D Frey, and Larry J Leifer. 2005. Engineering design thinking, teaching, and learning. *Journal of engineering education*, 94(1):103–120. 698
699
700
701

Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan, Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu, and Baobao Chang. 2025. **Omni-MATH: A universal olympiad level mathematic benchmark for large language models.** In *The Thirteenth International Conference on Learning Representations*. 702
703
704
705
706
707
708
709
710

Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu 711
712

713	Feng, Hanlin Zhao, and 1 others. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. <i>arXiv preprint arXiv:2406.12793</i> .	multimodal mathematical reasoner. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7126–7133.	769 770 771
716	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:2407.21783</i> .	Yong Lin, Shange Tang, Bohan Lyu, Jiayun Wu, Hongzhou Lin, Kaiyu Yang, Jia Li, Mengzhou Xia, Danqi Chen, Sanjeev Arora, and 1 others. 2025. Goedel-prover: A frontier model for open-source automated theorem proving. <i>arXiv preprint arXiv:2502.07640</i> .	772 773 774 775 776 777
721	Aryan Gulati, Brando Miranda, Eric Chen, Emily Xia, Kai Frondal, Bruno de Moraes Dumont, and Sanmi Koyejo. 2024. Putnam-axiom: A functional and static benchmark for measuring higher level mathematical reasoning. In <i>The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24</i> .	Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, and 1 others. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. <i>arXiv preprint arXiv:2405.04434</i> .	778 779 780 781 782 783
727	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts . In <i>The Twelfth International Conference on Learning Representations</i> .	784 785 786 787 788 789 790
733	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. <i>arXiv preprint arXiv:2103.03874</i> .	Pingchuan Ma, Tsun-Hsuan Wang, Minghao Guo, Zhiqing Sun, Joshua B Tenenbaum, Daniela Rus, Chuang Gan, and Wojciech Matusik. 2024. Llm and simulation as bilevel optimizers: A new paradigm to advance physical scientific discovery. <i>arXiv preprint arXiv:2405.09783</i> .	791 792 793 794 795 796
738	Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, and 1 others. 2025. Mathperturb: Benchmarking llms' math reasoning abilities against hard perturbations. <i>arXiv preprint arXiv:2502.06453</i> .	Yujun Mao, Yoon Kim, and Yilun Zhou. 2024. CHAMP: A competition-level dataset for fine-grained analyses of LLMs' mathematical reasoning capabilities . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 13256–13274, Bangkok, Thailand. Association for Computational Linguistics.	797 798 799 800 801 802
744	Yiming Huang, Zhenghao Lin, Xiao Liu, Yeyun Gong, Shuai Lu, Fangyu Lei, Yaobo Liang, Yelong Shen, Chen Lin, Nan Duan, and 1 others. 2023. Competition-level problems are effective llm evaluators. <i>arXiv preprint arXiv:2312.02143</i> .	Seyed Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2025. GSM-symbolic: Understanding the limitations of mathematical reasoning in large language models . In <i>The Thirteenth International Conference on Learning Representations</i> .	803 804 805 806 807 808
749	Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, and 1 others. 2024. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. 2024. Orca-math: Unlocking the potential of slms in grade school math. <i>arXiv preprint arXiv:2402.14830</i> .	809 810 811 812
755	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. Prometheus: Inducing fine-grained evaluation capability in language models . In <i>The Twelfth International Conference on Learning Representations</i> .	Nayantara Mudur, Hao Cui, Subhashini Venugopalan, Paul Raccuglia, Michael P Brenner, and Peter Norgaard. 2025. Feabench: Evaluating language models on multiphysics reasoning ability. <i>arXiv preprint arXiv:2504.06260</i> .	813 814 815 816 817
762	Ming Li, Jike Zhong, Tianle Chen, Yuxiang Lai, and Konstantinos Psounis. 2024. Eee-bench: A comprehensive multimodal electrical and electronics engineering benchmark. <i>arXiv preprint arXiv:2411.01492</i> .	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2080–2094.	818 819 820 821 822 823
767	Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. 2023. Unimath: A foundational and		

824	ZZ Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanxia Zhao, Liyue Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, and 1 others. 2025. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition. <i>arXiv preprint arXiv:2504.21801</i> .	881
825		882
826		883
827		884
828		885
829		886
830	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark . In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	887
831		888
832		889
833		
834		
835		
836	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, and 1 others. 2024. Deepseek-math: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	890
837		891
838		892
839		893
840		894
841		895
842	Saurabh Srivastava, Anto PV, Shashank Menon, Ajay Sukumar, Alan Philipose, Stevin Prince, Sooraj Thomas, and 1 others. 2024. Functional benchmarks for robust evaluation of reasoning performance, and the reasoning gap. <i>arXiv preprint arXiv:2402.19450</i> .	896
843		897
844		898
845		899
846		
847	Usman Syed, Ethan Light, Xingang Guo, Huan Zhang, Lianhui Qin, Yanfeng Ouyang, and Bin Hu. 2024. Benchmarking the capabilities of large language models in transportation system engineering: Accuracy, consistency, and reasoning behaviors. <i>arXiv preprint arXiv:2408.08302</i> .	900
848		901
849		902
850		903
851		904
852		905
853	Zhengyang Tang, Chenyu Huang, Xin Zheng, Shixi Hu, Zizhuo Wang, Dongdong Ge, and Benyou Wang. 2024. Orlm: Training large language models for optimization modeling. <i>arXiv preprint arXiv:2405.17743</i> .	906
854		907
855		908
856		909
857		
858	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	910
859		911
860		912
861		913
862		914
863		915
864	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	916
865		917
866		918
867		919
868		920
869		
870	Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with MATH-vision dataset . In <i>The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	921
871		922
872		923
873		924
874		925
875		
876	Xinlei Wang, Maike Feng, Jing Qiu, Jinjin Gu, and Junhua Zhao. 2024b. From news to forecast: Integrating event analysis in llm-based time series forecasting with reflection. <i>Advances in Neural Information Processing Systems</i> , 37:58118–58153.	
877		
878		
879		
880		
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024c. MMLU-pro: A more robust and challenging multi-task language understanding benchmark . In <i>The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	
	Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan J Prenger, and Animashree Anandkumar. 2023. Lendojo: Theorem proving with retrieval-augmented language models . <i>Advances in Neural Information Processing Systems</i> , 36:21573–21612.	
	Albert S Yue, Lovish Madaan, Ted Moskowitz, DJ Strouse, and Aaditya K Singh. 2024. Harp: A challenging human-annotated math reasoning benchmark . <i>arXiv preprint arXiv:2412.08819</i> .	
	Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, William Song, Tiffany Zhao, Pranav Raja, Charlotte Zhuang, Dylan Slack, and 1 others. 2024. A careful examination of large language model performance on grade school arithmetic. <i>Advances in Neural Information Processing Systems</i> , 37:46819–46836.	
	Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. 2022. minif2f: a cross-system benchmark for formal olympiad-level mathematics . In <i>International Conference on Learning Representations</i> .	
	Xiyuan Zhou, Huan Zhao, Yuheng Cheng, Yuji Cao, Gaoqi Liang, Guolong Liu, Wenxuan Liu, Yan Xu, and Junhua Zhao. 2024. Elecbench: a power dispatch evaluation benchmark for large language models . <i>arXiv preprint arXiv:2407.05365</i> .	

Appendix

Contents

A	The Use of Large Language Models	12
B	Ethical Considerations	12
C	Future Works	12
D	Dataset Construction	12
	D.1 Level 1 & Level 2 Extraction Process	12
	D.2 Level 3 Data Collection and Processing	13
	D.3 Version Variant Generation	14
E	Dataset URLs, License, and Hosting Plan	16
	E.1 Dataset Instance Metadata	16
F	Evaluation Details	17
	F.1 Level 1 & Level 2 Evaluation Details	17
	F.2 Level 3 Evaluation Details	17
	F.2.1 Rubric Construction	17
	F.2.2 Rubric-based Scoring and Human Calibration	18
	F.3 Level 3 Scoring Examples	18
G	Additional Analysis	22
	G.1 Level 3 Scoring Consistency Analysis	22
	G.2 Level 1 Analysis	22
	G.3 Level 2 Analysis	23
	G.4 Level 3 Analysis	23
	G.5 Subfield Performance Analysis	24

A The Use of Large Language Models

In this work, LLMs were used in three ways: (1) grammar checking and language polishing during paper writing, (2) generating controlled problem variants in the benchmark construction process, and (3) serving as both the models under evaluation and auxiliary judges for rubric-based scoring.

B Ethical Considerations

This work introduces a benchmark for evaluating large language models on engineering tasks. The problems are derived from publicly available benchmarks, academic competitions, and educational materials. For open-ended tasks, human participants voluntarily contributed reference solutions and evaluation scores using publicly available

rubric criteria, and personal information was collected only for inclusion in the acknowledgment section with explicit consent. The dataset does not contain sensitive data or enable harmful applications. EngiBench is designed as an evaluation framework for systematically analyzing and comparing model behaviors across diverse engineering task settings. The goal of EngiBench is to promote rigorous, transparent, and fair evaluation of language models in engineering contexts, and we affirm adherence to the ACL Code of Ethics, including principles of fairness, transparency, and research integrity.

C Future Works

While EngiBench establishes a strong foundation for evaluating LLMs on engineering problem-solving, several avenues remain for further development and expansion:

Scalability Across Engineering Domains. EngiBench currently covers three core engineering subfields—Systems & Control, Physical & Structural, and Chemical & Biological—which together span a wide range of disciplines such as Mechanical, Electrical, and Chemical/Biological Engineering. The benchmark framework is designed to be broadly applicable and adaptable across domains. In future work, we plan to expand the dataset by incorporating problems from additional engineering disciplines to further enhance data volume and subject diversity.

Multimodal Evaluation Extensions. Future versions of EngiBench will introduce a dedicated multimodal subset to evaluate models on tasks involving vision-language reasoning. This will enable systematic assessment of model performance in scenarios that demand visual interpretation alongside textual understanding.

Support for Long-Context Reasoning. We plan to extend the benchmark to include long-context engineering tasks by leveraging models with expanded context windows or hierarchical processing capabilities. This will allow for evaluation of more complex, information-rich tasks currently excluded due to input length limitations.

D Dataset Construction

D.1 Level 1 & Level 2 Extraction Process

To construct a high-quality and diverse dataset for Level 1 and Level 2, we systematically extract relevant tasks from a range of established public bench-

marks, including MMLU (Hendrycks et al., 2021), MATH (Hendrycks et al., 2021), GSM8k (Cobbe et al., 2021), Orca-Math (Mitra et al., 2024), HARP (Yue et al., 2024), Omni-MATH (Gao et al., 2025), Big-MATH (Albalak et al., 2025), and competition datasets such as cn_k12, Olympiads, AOPS forum, and AMC-AIME (Huang et al., 2023). In addition to these public sources, we also incorporate university-level engineering educational materials, including assignments, examinations, and instructor-provided teaching content, to further increase task diversity and real-world relevance.

To transform mathematical and logic-oriented problems into engineering-relevant evaluation tasks, we design a structured data processing pipeline that combines LLM-based analysis with human verification to ensure engineering relevance and classification accuracy. This pipeline ensures that all included problems align with real-world engineering semantics and reasoning demands, forming the basis for Level 1 and Level 2 in EngiBench.

The processing pipeline consists of the following steps:

- 1. Engineering Relevance Filtering:** Each problem is evaluated for its applicability to engineering scenarios. Problems lacking domain relevance are excluded to maintain the technical integrity of the benchmark. The prompt used to determine whether a problem pertains to engineering is as follows:

```
1 """Determine if ORIGINAL problem
2 can be solved with ONLY
3 mathematical knowledge (NO
4 engineering background):
5 - False if requires any domain-
6 specific knowledge
7 - True if solvable through pure
8 mathematical calculations"""
```

- 2. Discipline and Subfield Classification:** Relevant problems are first assigned to a specific engineering discipline (e.g., Electrical, Civil, Mechanical), and then grouped into one of EngiBench’s three high-level analytical subfields: Systems & Control, Physical & Structural, or Chemical & Biological. The prompt used for assigning a problem to a specific engineering discipline is as follows:

```
1 """If yes, which engineering
2 category? (Chemical/
3 Bioengineering/Geotechnical/
4 Energy/Nuclear/Aerospace/
5 Automotive/Biomedical/Civil/
```

```
Control/Electrical/Industrial/
Mechanical/Ocean/Environmental/
Other) (Please try to avoid
Other)
2 If not an engineering problem,
3 return "N/A"."""
```

- 3. Difficulty Level Assignment:** Based on the complexity of the required reasoning process, tasks are categorized into Level 1 or Level 2. Level 1 includes basic knowledge recall and single-step computation, while Level 2 involves multi-step inference, contextual understanding, and integration of structured constraints. The prompt used for classifying the difficulty level of a problem is as follows:

```
1 """Difficulty level? (Level 1/
2 Level 2) (Please try to avoid
3 unknown):
4 - Level 1: The problem can be
5 solved by a direct retrieval of
6 information or by directly
7 substituting values into a known
8 formula i.e., the shortest
9 possible solution path. No
10 chaining of intermediate steps
11 is required. (Example: Using Ohm
12 's Law,  $V = IR$ , to directly
13 compute voltage when given
14 current and resistance.)
15 - Level 2: The problem requires
16 multi-step reasoning meaning
17 that it involves chaining
18 together several logical
19 deductions, intermediate
20 calculations, or systematic
21 strategies beyond a single
22 direct formula application. (
23 Example: Analyzing a circuit to
24 compute total resistance by
25 first calculating individual
26 branch resistances and then
27 combining them.)"""
```

D.2 Level 3 Data Collection and Processing

To construct the Level 3 dataset in EngiBench, we focus on real-world, open-ended engineering tasks sourced from major mathematical modeling competitions. Specifically, we collect problems from publicly accessible archives of contests such as the China Undergraduate Mathematical Contest in Modeling (CUMCM), the Mathematical Contest in Modeling / Interdisciplinary Contest in Modeling (MCM/ICM), and the Asia and Pacific Mathematical Contest in Modeling (APMCM), covering the years 2010 to 2024.

To ensure domain relevance and evaluation consistency, we apply strict filtering criteria. We retain

only problems with clear engineering context and official scoring rubrics, and exclude those that depend heavily on complex diagrams or large external tables requiring multimodal input.

We standardize the selected problems using a structured pipeline that combines LLM-based processing with human oversight. This ensures language clarity, formatting consistency, and reduced risk of data contamination. The pipeline includes the following steps:

- 1. Language Normalization:** Non-English problems are translated into fluent English using machine translation, while preserving the original engineering semantics.
- 2. Expression Rewriting:** To minimize potential overlap with pretraining data, each problem is paraphrased by the LLM using diverse sentence structures and reasoning styles. While surface expressions are significantly altered, the core logic, numerical values, and solution paths remain unchanged. This step produces the *perturbed version* of each task, which is used to evaluate model robustness to superficial input variations.
- 3. Multimodal Simplification:** For problems containing simple figures or tables, we extract and describe the essential information using plain text or \LaTeX -formatted representations to support uniform text-based evaluation.

LLM Prompt Template: The following instruction prompt is used to guide the LLM in modifying each problem:

```
1 """Assuming you are a question
2 expert, please translate this
3 question into English. And while
4 ensuring that the meaning of the
5 question remains unchanged (
6 preserving all logic, values, and
7 the type of reasoning required),
8 change the way the question is
9 expressed by rewriting it in a way
10 that is radically different from
11 your regular logical structure,
12 simulating the randomness of manual
13 rewriting by human experts, and
14 using as many sentence variations as
15 possible. If there is a table,
16 please convert it into a table form
17 using LaTeX. For simple pictures,
18 please describe them directly. The
19 question is required to be converted
20 into is in str format."""
```

To ensure the technical rigor and domain consistency of the Level 3 dataset, the entire generation and transformation process was closely supervised and iteratively revised by doctoral-level professionals with extensive expertise in engineering and mathematical modeling. These experts reviewed both the selection of source problems and the outputs produced by the language model, verifying that each task preserved the original problem’s intent, accurately reflected real-world engineering reasoning, and met the standards expected in academic and professional modeling contexts.

The details of how the original contest scoring standards were mapped into EngiBench’s formal scoring rubrics are described in the later subsection (see Section F.2).

D.3 Version Variant Generation

To assess model robustness and isolate specific reasoning limitations, we generate three structured variants for each Level 1 and Level 2 problem: *Perturbed*, *Knowledge-Enhanced*, and *Math Abstraction*. These variants are created through LLM prompting, with manually verified outputs to ensure alignment with the original problem logic and correctness. Below, we describe the purpose and generation criteria for each variant, accompanied by illustrative prompts.

- **Perturbed Variant.** This variant alters the surface form of the original problem—either through numerical or linguistic changes—while preserving its core logic and computational requirements. The purpose is to test whether model performance stems from true reasoning ability or superficial pattern matching. A rewriting suitability code (0–3) guides the type of modification to apply. The prompt used to generate the perturbed version and related content is as follows:

```
1 """
2 1. Rewriting Suitability: Determine
3 the type (0-3):
4 - 0: Non-rewritable (use only
5 when necessary)
6 - 1: Modify expressions only
7 - 2: Modify numerical values only
8 - 3: Modify both expressions and
9 numerical values
10 // Note: All rewrites must
11 maintain the original problem
12 logic, engineering context, and
13 reasoning/computational
14 requirements
```

```

1234 9 2. Rewritten Problem: Rewrite the
1235     problem according to the type of
1236     rewriting suitability above.
1237     Make the answer as difficult as
1238     possible while ensuring that the
1239     answer is correct. (Please
1240     rewrite the problem in a way
1241     that is radically different from
1242     your regular logical structure
1243     by: (1) avoiding common
1244     reasoning patterns in your model
1245     , (2) simulating human expert
1246     manual rewriting randomness, and
1247     (3) using maximum sentence
1248     variation.)
1249     - If 0, return original problem
1250     unchanged
1251     - If 1, modify expressions only
1252     - If 2, modify numerical values
1253     only
1254     - If 3, modify both expressions
1255     and values
1256
1257 15 3. Rewritten Solution Process:
1258     Provide step-by-step explanation
1259     including all reasoning,
1260     calculations and logic. Clearly
1261     state if answer can be obtained
1262     directly through formula
1263     substitution (shortest solution
1264     path without intermediate steps)
1265     .
1266
1267 17 4. Rewritten Answer: Provide correct
1268     answer for rewritten problem (
1269     only types 2/3 may change)"""
1270

```

1271 • **Knowledge-enhanced Variant.** In this version, relevant domain knowledge—such as formulas, constants, and conversions—is explicitly provided before the original question. This allows us to evaluate whether performance deficits are due to missing knowledge or failures in application. The question itself is unchanged to isolate the impact of added context. The prompt used to generate the knowledge-enhanced version is as follows:

```

1281 1 """Knowledge-Enhanced Version:
1282 2 WARNING: Make sure the final
1283     numerical answer to the
1284     converted mathematical problem
1285     is exactly the same as the
1286     original problem.
1287
1288 3
1289 4 Given:
1290 5 - List all relevant formulas or
1291     principles (e.g., Ohm's Law:  $V = I * R$ )
1292 6 - Include physical constants with
1293     values if they are involved (e.g
1294     .,  $g = 9.8 \text{ m/s}^2$ )
1295 7 - Specify unit conversions if
1296     applicable (e.g.,  $1 \text{ kWh} = 3.6 * 10^6 \text{ J}$ )
1297

```

```

8 - State any assumptions or ideal
1298     conditions if necessary (e.g.,
1299     assume no heat loss)
1300
1301 9
1302 10 Problem:
1303 11 Repeat the original question exactly
1304     as stated
1305
1306 12
1307 13 Example:
1308 14 Original: "Calculate voltage across
1309     5 Ohm resistor with 2 A current"
1310 15 Enhanced:
1311 16 "Given:
1312 17 - Ohm's Law:  $V = I * R$ 
1313 18 - Problem: Calculate voltage across
1314     5 Ohm resistor with 2 A current"
1315 19 """"

```

1316 • **Math Abstraction Variant.** This version reformulates the original engineering problem into a purely mathematical format by removing all domain-specific context. Variables and operations are explicitly defined to preserve the exact calculation logic. This allows us to isolate whether reasoning failure arises from contextual understanding or mathematical ability. The prompt used to generate the math abstraction version is as follows:

```

1326 1 """"Rewrite the given problem into a
1327     purely mathematical version by:
1328
1329 2
1330 3 a. Remove all domain-specific
1331     context (e.g., chemistry,
1332     physics, economics).
1333 4 b. Keep only numbers, variables, and
1334     math operations.
1335 5 c. If domain-specific knowledge is
1336     required (e.g., reaction ratio,
1337     atomic mass), extract only the
1338     final numerical ratio or
1339     constant and include it directly
1340     .
1341 6 d. Maintain the exact calculation
1342     logic and final answer.
1343 7 e. Use structured symbolic language
1344     in a compact form:
1345 8 - Introduce variables explicitly (e.
1346     g., "Let  $x = 2$  and  $y = 3.$ ")
1347 9 - Define the calculation clearly (e.
1348     g., "Total  $z = \min(x, y) * 2.$ ")
1349 10 - End with "Find the result."
1350 11
1351 12 WARNING: Make sure the final
1352     numerical answer to the
1353     converted mathematical problem
1354     is exactly the same as the
1355     original problem.
1356 13
1357 14 Examples:
1358 15 Original: "In the reaction:  $\text{Cl}_2 + \text{H}_2$ 
1359     ->  $2\text{HCl}$ , 1 mole of  $\text{Cl}_2$  reacts
1360     with 2 moles of  $\text{H}_2$ . How many
1361     moles of  $\text{HCl}$  can be formed?"

```

```

17 converted_problem: "Let x = 1 and y
1362 = 2. They react in the ratio x :
1363 y : z = 1 : 1 : 2. Total
1364 product z = min(x, y) * 2. Find
1365 the result."
1366
1367
18
1368 19 Original: "A 2m wide platform sinks
1369 0.01m under 60kg. Estimate its
1370 length assuming water density =
1371 1000 kg/m^3."
1372 20 converted_problem: "Let x = 60 / (2
1373 * 0.01 * 1000). Find the result
1374 ." ""
1375 21

```

E Dataset URLs, License, and Hosting Plan

EngiBench is released for research and evaluation purposes only. All third-party artifacts are used in accordance with their original licenses. The released benchmark does not redistribute restricted original content, and commercial use of the benchmark is not permitted.

E.1 Dataset Instance Metadata

For the EngiBench dataset, each instance corresponds to an engineering task and is stored in a structured format. Instances are categorized according to task difficulty (Level 1, 2, or 3) and are constructed with multiple versions to enable fine-grained evaluation of different capabilities. The metadata fields for each level are described below:

Level 1 and Level 2 Each row in the Level 1 & 2 dataset corresponds to a closed-form or structured engineering problem, and includes the following fields:

- **problem** – Original natural language problem statement.
- **answer** – Ground truth answer to the original problem.
- **subfield** – Engineering subfield to which the problem belongs (e.g., Systems & Control).
- **category** – Topic-specific classification within the subfield (e.g., Thermodynamics).
- **difficulty** – Either Level 1 (Foundational Knowledge Retrieval) or Level 2 (Contextual Reasoning).
- **converted_problem** – Abstract mathematical formulation of the problem.

- **converted_problem_llm_answer** – LLM-generated response to the converted problem. 1409 1410
 - **knowledge_enhanced_problem** – Problem reformulated with explicit formulas and domain definitions. 1411 1412 1413
 - **rewritten_problem** – Semantically or numerically perturbed variant of the original problem. 1414 1415 1416
 - **rewritten_answer** – Answer to the rewritten problem. 1417 1418
 - **rewritten_converted_problem** – Mathematical abstraction of the rewritten problem. 1419 1420
 - **rewritten_converted_problem_llm_answer** – LLM response to the rewritten converted problem. 1421 1422 1423
 - **rewritten_knowledge_enhanced_problem** – Knowledge-enhanced version of the rewritten problem. 1424 1425 1426
- Level 3** Each Level 3 instance represents an open-ended modeling task and includes both the problem prompt and a rubric-based evaluation across multiple capability dimensions: 1427 1428 1429 1430
- **question** – English translation of the open-ended modeling task. 1431 1432
 - **question_modified** – Semantically perturbed variant of the task. 1433 1434
 - **source_detail** – Source of the modeling task (e.g., MCM, coursework). 1435 1436
 - **official_scoring_standard** – English translation of rubric criteria. 1437 1438
 - **subfield** – Engineering subfield of the task. 1439
 - **category** – Domain or topic under which the task is categorized. 1440 1441
 - **information_extraction_score** – Score for identifying relevant variables and constraints. 1442 1443
 - **multi_objective_decision_score** – Score for resolving trade-offs across objectives. 1444 1445
 - **uncertainty_handling_score** – Score for reasoning under ambiguity or variable inputs. 1446 1447
 - **domain_specific_reasoning_score** – Score for applying engineering-specific logic and formulas. 1448 1449 1450

F Evaluation Details

F.1 Level 1 & Level 2 Evaluation Details

Level 1 and Level 2 tasks consist of well-defined problems with clearly defined and unique solutions. We therefore adopt a *binary scoring* scheme, where each model-generated answer is compared against a reference answer and marked as either correct (1) or incorrect (0). Overall performance is reported in terms of accuracy.

Evaluation is conducted through an automated comparison procedure. To handle diverse numerical formats, units, and equivalent expressions, we design a standardized evaluation prompt, which is independently executed by GPT-4.1 and Gemini 2.5 Flash. For cases where the two evaluators produce inconsistent judgments, manual verification is performed to determine the final decision. The evaluator determines whether a generated answer matches the reference answer based on mathematical correctness, unit validity, and logical consistency. For numerical questions, a tolerance of $\pm 2\%$ is allowed to account for rounding effects in multi-step calculations. The evaluator is instructed to output only a Boolean decision (“True” or “False”) to ensure consistent and reproducible scoring.

To verify evaluation consistency and reliability, we perform multi-model cross-checking and human spot checks. Specifically, all Level 1–2 responses are independently evaluated by GPT-4.1 and Gemini 2.5 Flash, and their results are compared. In addition, we randomly sample 300 problems for manual verification. On this subset, GPT-4.1 achieves an evaluation accuracy of 98.67%, while Gemini 2.5 Flash achieves 98.33%. These results demonstrate consistent evaluation behavior and show that the automated procedure closely aligns with deterministic answer matching for Level 1 and Level 2 tasks.

```
1 """Please analyze these two answers
1490   carefully:
1491   2 Generated Answer: {generated_answer}
1492   3 Standard Answer: {correct_answer}
1493   4
1494   5 Follow these rules for comparison:
1495   6 1. For calculation-focused problems:
1496   7   - If the numerical values match,
1497   8     consider it correct even if units
1498   9     are missing
1499   8   - Focus on the mathematical reasoning
1500   9     and final numerical result
1501   9   - Check if the core calculation steps
1502   10  are correct
1503   10  - For complex calculations, allow 2
1504   1505  % tolerance in the final numerical
```

```
11 result
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
3335
3336
3337
3338
3339
3340
3341
3342
3343
3344
3345
3346
3347
3348
3349
3350
3351
3352
3353
3354
3355
3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455
3456
3457
3458
3459
3460
3461
3462
3463
3464
3465
3466
3467
3468
3469
3470
3471
3472
3473
3474
3475
3476
3477
3478
3479
3480
3481
3482
3483
3484
3485
3486
3487
3488
3489
3490
3491
3492
3493
3494
3495
3496
3497
3498
3499
3500
3501
3502
3503
3504
3505
3506
3507
3508
3509
3510
3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563
3564
3565
3566
3567
3568
3569
3570
3571
3572
3573
3574
3575
3576
3577
3578
3579
3580
3581
3582
3583
3584
3585
3586
3587
3588
3589
3590
3591
3592
3593
3594
3595
3596
3597
3598
3599
3600
3601
3602
3603
3604
3605
3606
3607
3608
3609
3610
3611
3612
3613
3614
3615
3616
3617
3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671
3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725
3726
3727
3728
3729
3730
3731
3732
3733
3734
3735
3736
3737
3738
3739
3740
3741
3742
3743
3744
3745
3746
3747
3748
3749
3750
3751
3752
3753
3754
3755
3756
3757
3758
3759
3760
3761
3762
3763
3764
3765
3766
3767
3768
3769
3770
3771
3772
3773
3774
3775
3776
3777
3778
3779
3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
38
```

```

1565 2
1566 3 For each capability that is covered,
1567 4 provide a scoring rubric in the
1568 5 following format:
1569 6
1570 7 Problem [(Problem ID)]:
1571 8 redundant_information_filtering_score:
1572 9 (1)(2)...
1573 10 multi_objective_tradeoff_score: (1)(2)
1574 11 ...
1575 12 uncertainty_handling_score: (1)(2)...
1576 13 deep_knowledge_integration_score: (1)(2)
1577 14 ...
1578 15
1579 16 Notes: Each capability has a total
1580 17 possible score of 10 points. In
1581 18 other words, the total score for
1582 19 each listed capability should sum to
1583 20 10 points. Capabilities that are
1584 21 not covered in this problem receive
1585 22 0 points. The rubric should further
1586 23 specify, under each capability, the
1587 24 different score levels (e.g., 1
1588 25 point, 2 points, 3 points, etc.) and
1589 26 the corresponding specific
1590 27 behaviors or response
1591 28 characteristics associated with each
1592 29 level.
1593 30
1594 31 Please read the problem and rubric
1595 32 carefully and provide a capability-
1596 33 based evaluation rubric for how this
1597 34 problem assesses the output of
1598 35 large language models.""""

```

F.2.2 Rubric-based Scoring and Human Calibration

The finalized rubrics are applied to evaluate model-generated responses for Level 3 tasks. We implement an automated LLM-based scoring pipeline that assesses solution quality along multiple capability dimensions defined by the rubrics. Specifically, scores are independently produced by GPT-4.1 and Gemini 2.5 Flash, and the final score is obtained by averaging the two to reduce evaluator-specific variability.

To ensure the reliability of the reported results, LLM-generated scores are reviewed and calibrated by annotators with engineering backgrounds. The main results in this paper report calibrated scores. We note that fully automated LLM-based scoring already provides a strong and practical reference, as further supported by the consistency and validity analysis in Appendix G.1.

The prompt used to evaluate the generated answer against the rubric is as follows:

```

1620 1 f """
1621 2 You are a professional modeling
1622 3 competition judge with extensive
1623 4 experience in evaluating
1624 5 mathematical and engineering models.

```

```

Please conduct a rigorous
evaluation of the following answer
based on the provided criteria.

Answer to evaluate:
{answer}

Evaluation Criteria:
{score_criteria}

Please evaluate strictly according
to the criteria and provide your
assessment in the following JSON
format:
{{
  "score": <score between 0-10,
can use decimal points for precision
>,
  "reason": "Detailed evaluation
breakdown:\n
          1. [Specific criterion
] - [sub-score] points: [
justification]\n
          2. [Specific criterion
] - [sub-score] points: [
justification]\n
          3. [Specific criterion
] - [sub-score] points: [
justification]\n
          Final score: [total]
points"
}}

Note:
- Break down your scoring into
specific components
- Provide clear justification for
each sub-score
- Be objective and consistent in
your evaluation
- Consider both the technical
accuracy and the methodology
"""

```

F.3 Level 3 Scoring Examples

As results shown in section 4.2.3, the answers of LLMs to open-ended tasks show significant differences in four dimensions of information extraction, multi-objective decision making, uncertainty handling and domain-specific reasoning. Figure 7 preliminarily presents two scoring segments, 3 points and 8 points, for the evaluation of models' answers. To demonstrate the response performance of different segments more clearly and intuitively, we provide the following examples with more Level 3 scoring details:

1. **Full Mark (Avg. Score: 9.475):** The problem requires optimizing Hu sheep farm pen utilization under stochastic conditions (conception rates, gestation periods, litter sizes) while adhering to strict capacity constraints and cohabitation rules. The solution must minimize

1686 expected losses from idle pens (1 unit/day)
 1687 or shortages (3 units/day) through dynamic
 1688 scheduling and statistical validation.

1689 • **Information Extraction (10/10):**

1690 Exclusion of Deterministic Assumptions
 1691 (5/5): Section 1 (System Overview)
 1692 clarifies all critical parameters mod-
 1693 eled as random variables (e.g., “ $X_c \sim$
 1694 $\text{Binomial}(N_m, 0.85)$: Number of suc-
 1695 cessful conceptions; $G \sim U[147, 150]$:
 1696 Gestation days; $L_s \sim \text{Poisson}(\lambda = 2.2)$:
 1697 Liveborn lambs per ewe, with 3% mor-
 1698 tality ($L_a = L_s \cdot 0.97$); $L_d \sim U[35, 45]$:
 1699 Lactation days”). Section 3A (Scenario
 1700 Generation) replaces fixed values with
 1701 dynamic sampling (e.g., “For each sce-
 1702 nario, sample: - Which ewes conceive
 1703 (Bernoulli, 85%) - Their gestation (G) -
 1704 Number of lambs (L_s), apply mortality -
 1705 Lactation length (L_d)”). Section 6B (Ro-
 1706 bust Planning) makes flexible scheduling
 1707 responsive to stochastic outcomes (e.g.,
 1708 “Adjust mating/rest period within allowed
 1709 windows to shift animal flows.”).

1710 Identification of Valid Uncertainty
 1711 Parameters (5/5): Section 1 clarifies
 1712 explicit distributions for all uncertainties
 1713 (e.g., “ $X_c \sim \text{Binomial}(N_m, 0.85)$...
 1714 $G \sim U[147, 150]$... $L_s \sim$
 1715 $\text{Poisson}(2.2)$... $L_d \sim U[35, 45]$ ”).
 1716 Section 3A ensures consistent appli-
 1717 cation in scenario generation (e.g.,
 1718 “Sample conception (Bernoulli), ges-
 1719 tation (G), litter size (L_s), lactation
 1720 (L_d)”). Section 5 (Loss Function) offers
 1721 loss calculation integrating stochastic
 1722 inputs (e.g., “ $\mathbb{E}_{\text{scenario}} [\sum_t [I_t + 3S_t]]$ ”).

1723 • **Multi-objective Decision making**
 1724 **(9.2/10):**

1725 Minimized Expected Loss & Output
 1726 Maximization (4.5/5): Section 5 (Loss
 1727 Function) contains rigorous mathemat-
 1728 ical formulation balancing idle (1 unit)
 1729 vs. shortage (3 unit) costs (e.g., “Objec-
 1730 tive: $\min \mathbb{E}_{\text{scenario}} [\sum_t [I_t + 3S_t]]$ $I_t =$
 1731 Idle pens, $S_t = \text{Shortages}$ ”). Section 7B
 1732 (Robust Planning) includes statistical val-
 1733 idation of tradeoffs (e.g., “Monte Carlo
 1734 over Scenarios: Simulate losses across
 1735 all scenarios for each candidate policy.”)
 1736 Section 8 (Results Table) applies quanti-

tative comparison of policies.

Lactation Flexibility & Fattening
 Tradeoffs (4.7/5): Section 1 (System
 Overview) makes explicit dynamic
 linkage between lactation and fattening
 (e.g., “ $L_d \sim U[35, 45]$: Lactation days
 $\rightarrow F_d = 210 + 2 \cdot (40 - L_d)$: Fattening
 days”). Section 6B (Robust Planning)
 considers operational use of flexibility
 to smooth demand (e.g., “Adjust rest
 periods to align cohorts, minimizing
 ‘loner pens.’”). Section 3A (Scenario
 Generation) has stochastic integration of
 tradeoff (e.g., “Sample lactation length
 (L_d), impact on fattening (F_d)”).

• **Uncertainty Handling (9.2/10):**

Stochastic Process Models (4/4): Section
 1 (System Overview) specifies explicit
 distributions for all stochastic parame-
 ters (e.g., “ $X_c \sim \text{Binomial}(N_m, 0.85)$,
 $G \sim U[147, 150]$, $L_s \sim \text{Poisson}(2.2)$,
 $L_d \sim U[35, 45]$ ”). Section 3A (Sce-
 nario Generation) implements full Monte
 Carlo (e.g., “Generate 1000 scenarios...
 sample conception (Bernoulli), gestation
 (G), litter size (L_s), lactation (L_d)”).
 Section 7B (Robust Planning) includes
 statistical validation of stochastic out-
 comes (e.g., “For each candidate policy,
 simulate losses across all scenarios.”).

Dynamic Adjustment Strategies (2.7/3):
 Section 1 (Fattening Calculation) estab-
 lishes mechanistic linkage of lactation-
 fattening tradeoff (e.g., “ $F_d = 210 + 2 \cdot$
 $(40 - L_d)$: Fattening days adjusted by lac-
 tation.”). Section 6B (Robust Planning)
 makes adaptive scheduling but lacks two-
 way feedback (e.g., “Adjust rest peri-
 ods to align cohorts... weekly rolling
 re-optimization.”).

Contingency Sets (2.5/3): Section 2 (Co-
 habitation Rules) contains hard-coded
 tolerance for uncertainty (e.g., “Group
 into largest feasible penfuls within 7-day
 windows.”). Section 8 (Statistical Assess-
 ment) analyzes multi-scenario sensitivity
 (e.g., “Tabulate average loss, shortage
 probability, and max pen use.”).

• **Domain-specific Reasoning (9.5/10):**

Integration of Empirical Rules (4/4):
 Section 2 (Cohabitation Rules) adds

hard-codes industry constraints into algorithms (e.g., “7-day tolerance window for nursing ewes, lambs, and resting ewes... Group into largest feasible penfuls (14 fattening lambs/pen, 6 nursing ewes/pen).”). Section 1 (System Overview) uses embeds empirical flexibility ranges as distributions (e.g., “ $L_d \sim U[35, 45]$: Lactation days... $R \sim U[18, 22]$: Adjustable rest period.”) Section 6B (Robust Planning) operationalizes flexible rest rules (e.g., “Extend rest periods to align cohorts if pens would otherwise idle.”).

Expected Loss Functions (3/3): Section 5 (Loss Function) has rigorous probabilistic loss aggregation (e.g., “ $\min \mathbb{E}_{\text{scenario}} [\sum_t [I_t + 3S_t]]$ ”, $I_t = \max(P_{\text{avail}} - P_{\text{req}}(t), 0)$, $S_t = \max(P_{\text{req}}(t) - P_{\text{avail}}, 0)$.”). Section 8 (Results Table) quantifies loss distribution across scenarios. Section 3B (State Evolution) links stochastic occupancy to loss calculation (e.g., “For each day t : Compute $P_{\text{req}}(t)$ from sampled cohorts.”).

Stochastic Optimization Algorithms (2.5/3): Section 7B (Robust Planning) applies sample average approximation (SAA) method (e.g., “Monte Carlo simulation over 1000 scenarios to evaluate policies.”). Section 6A (Rolling Horizon) uses heuristic dynamic programming (e.g., “Re-optimize mating batches weekly to maximize cohabitation.”).

2. **5 points (Avg. Score: 5.375):** The problem involves modeling a team coordination exercise (“Unity Drum”) where 8 members control a drum’s tilt by pulling ropes to bounce a ball. Key tasks include: 1. Calculating the drum’s tilt angle at $t=0.1s$ based on force/timing inputs (Table 1), accounting for initial 11cm displacement. 2. Ensuring physics-based accuracy in torque, angular acceleration, and geometric relationships.

• **Information Extraction (7.5/10):**

Error Source Analysis (5/6): Explicit Recognition: Timing errors-“Some members may apply force slightly before others” (Algorithm section); strength

variation-“Members likely have different strengths” (Considerations). Partial Implementation: Timing logic in code (if $\text{timing}[i] \leq 0.1$) is noted but lacks vector-time coupling; force scaling ($\text{effective_force} = \frac{\text{force}(\text{member_id}-1)}{10}$) is arbitrary.

Physical Model Simplification (2.5/4): Justified Simplifications: “Ignores damping for short-duration calculation” (Considerations); Drum as uniform cylinder ($I = 0.5 \cdot \text{drum_mass} \cdot r^2$). Over-Simplifications: Fixed torque angle ($\sin(\frac{\pi}{2})$) ignores vector geometry; rope tautness assumption (“If the drum tilts too far, ropes could slack”) not modeled.

• **Multi-objective Decision making (6.5/10):**

Tilt Angle and Force Relationship (4.5/6): Physics Foundation: Correctly derives torque ($\tau = r \cdot F \cdot \sin(\theta)$), inertia ($I = 0.5 \cdot m \cdot r^2$), and angular kinematics ($\theta = \theta_0 + \frac{1}{2}\alpha t^2$); maps rope geometry ($\text{angle_radians} = (\text{member_id} - 1) \cdot (\frac{2\pi}{8})$). Implementation Gaps: Timing logic (if $\text{timing}[i] \leq 0.1$) is crude; forces are binary (on/off) rather than time-interpolated; no optimization for tilt minimization (e.g., predictive control or force balancing).

Computational Efficiency (2/4): Basic Looping-iterates over 8 members with $O(1)$ operations per member (e.g., $\text{torque} = \text{drum_radius} \cdot \text{force} \cdot \sin(\frac{\pi}{2})$). No Advanced Techniques-lacks vectorization, memoization, or scalability for larger teams.

• **Uncertainty Handling (2/10):**

Error Propagation Analysis (2/4): Acknowledgment Only: Mentions “members likely have different strengths and reaction times” (Considerations); suggests “extended to simulate more realistic distributions” but provides no math or implementation. No Quantification: Lacks sensitivity analysis or error bounds on tilt angle.

Numerical Simulation Estimation (0/4): No Monte Carlo: Code calculates tilt for fixed inputs only (force_data); no randomization of force/timing or statistical

1889	output (mean/variance).	leaved Logs Without Justification: The	1939
1890	Methodological Clarity (N/A): Physics	primary and secondary transmission logs	1940
1891	steps are clear but irrelevant to uncer-	(Tables 1-2) are interleaved in the solu-	1941
1892	tainty scoring.	tion ("Round 1: Primary 1→2; Round 1:	1942
1893	• Domain-specific Reasoning(5.5/10):	Secondary 1→1a"), but no protocol en-	1943
1894	3D Mechanics Modeling (2.5/6): 2D	sures collision avoidance (e.g., TDMA,	1944
1895	Limitation: Explicitly states "our coord-	priority scheduling). Unverified Simul-	1945
1896	inate system will be planar (X and Y	taneity Assumption: The answer states	1946
1897	only)" (Key Equations); torque calcula-	"Simultaneous reception allowed during	1947
1898	tion ($\tau = r \cdot F \cdot \sin(\theta)$) ignores out-of-	transmission" (Step 1) but doesn't prove	1948
1899	plane forces. Partial Physics: Correctly	this suffices for concurrent primary/sec-	1949
1900	models drum as cylinder ($I = 0.5 \cdot m \cdot r^2$)	ondary transmissions under the 8-minute	1950
1901	but lacks 3D rotation dynamics.	constraint.	1951
1902	Model-Based Optimization Strategy	• Multi-objective Decision making	1952
1903	(3/4): Suggestions Without Implemen-	(2/10):	1953
1904	tation: Proposes "damping term propor-	3D Parameter Optimization (0/6):	1954
1905	tional to angular velocity" (Considera-	Single-Parameter Focus: The answer	1955
1906	tions); mentions "member variation" but	only optimizes for N_max	1956
1907	no adaptive control (e.g., PID for tilt cor-	(" $N(N-1)/28 \rightarrow N_{max} = 4$ ", Step	1957
1908	rection).	2) but ignores joint optimization of	1958
1909	3. 1 point (Avg. Score: 1.25): The problem in-	capability (no analysis of 158-character	1959
1910	volves coordinating multiple meteorological	message limits or segment splitting	1960
1911	units (each with 1 primary and 2 secondary	efficiency), reliability (no adjustment	1961
1912	stations) to ensure reliable hourly weather	for secondary station 80% success rate	1962
1913	data collection and full data sharing under	such as no retransmission strategy) and	1963
1914	strict communication constraints. Key chal-	time (assumes 8 minutes suffice without	1964
1915	lenges include managing transmission reliabil-	validating secondary transmission	1965
1916	ity (80% for secondaries, 100% for primaries),	overhead). Missed Pareto Frontier: Fails	1966
1917	message capacity limits, and achieving 97%	to explore tradeoffs (e.g., "Could N=5	1967
1918	success probability within 8 minutes for pri-	work if secondary transmissions are	1968
1919	mary data exchange. The goal is to determine	reduced?").	1969
1920	the maximum number of units (Nmax), de-	Resource Allocation Strategy (2/4):	1970
1921	sign transmission schemes, and compute per-	Equal Bandwidth Only: Primary stations	1971
1922	formance metrics.	follow a round-robin schedule ("1→2,	1972
1923	• Information Extraction (2/10): High-	1→3, 1→4, 2→3, ...", Table 1), and	1973
1924	Probability Constraint Processing (0/5):	secondaries transmit uniformly ("1→1a,	1974
1925	Failure to Address Probabilistic Guar-	1→1b, 2→2a, ...", Table 2). No Prioriti-	1975
1926	antee: The answer calculates secondary	zation: Critical objectives (e.g., ensur-	1976
1927	transmission success as "expected num-	ing 97% success) aren't prioritized in	1977
1928	ber of reports received... is $4 \times 0.8 = 3.2$ "	scheduling.	1978
1929	(Step 4) but never models retransmis-	• Uncertainty Handling (0/10):	1979
1930	sions or redundancy to achieve 97% suc-	High-Order Probability Events (0/6):	1980
1931	cess. The assumption of direct success	No Threshold Calculation: The answer	1981
1932	ignores the problem's explicit probabil-	states secondary stations have an	1982
1933	ity requirement. Missing Critical Logic:	"80% transmission/reception success rate"	1983
1934	No discussion of how to compensate for	(Step 1) but never computes the prob-	1984
1935	the 20% failure rate (e.g., retrying failed	ability of achieving 97% success (e.g.,	1985
1936	transmissions, acknowledgments, or er-	via binomial distribution for multiple re-	1986
1937	ror correction).	tries). Misleading Metric: The "mean	1987
1938	Time Window Isolation (2/5): Inter-	secondary reports received per primary	1988
		station (3.2)" (Step 4) is irrelevant to the	1989

cumulative success probability requirement.

Asymmetric Loss (0/4): No Cost Analysis: The solution ignores idle time cost (unused transmission slots due to failures) and rental loss (penalties for delayed data delivery implied by "critical rescue operations").

• **Domain-specific Reasoning (1/10):**

Mixed-Integer Programming (0/5): No Optimization Model: The answer derives $N_{\max} = 4$ via a simple inequality (" $\frac{N(N-1)}{2} \leq 8$ ", Step 2) but lacks an objective function (e.g., "maximize N while meeting time/reliability constraints"), and omits integer constraints (N must be discrete) or linear relaxation techniques. Ad-Hoc Calculation: No use of MINLP (Mixed-Integer Nonlinear Programming) to jointly optimize N, transmission scheduling, and reliability. Fault-Tolerant Protocol Design (1/5): Basic Segmentation: Mentions "reports can split into two 50-character segments" (Step 1) but no dual verification (never states if segments are sent redundantly to different primaries) and no formal protocol (assumes secondary stations report to all primaries without fault recovery like checksums, ACKs).

G Additional Analysis

G.1 Level 3 Scoring Consistency Analysis

To further examine the reliability of our evaluation protocol, we compare three scoring variants for Level 3 tasks: human-calibrated scoring (reported as the main results), fully automated LLM-only scoring, and fully manual human scoring (human-only).

Table 2 summarizes the average scores under these three scoring settings for both original and perturbed tasks. Across models and task settings, LLM-only scores are generally close to human-calibrated scores, with differences typically within a small margin. This indicates that the proposed rubric enables reliable automated evaluation, while human calibration mainly serves to correct a limited number of edge cases and ensure maximum rigor in the reported results.

Regarding the validation of scoring consistency, the relative ordering of models remains largely con-

sistent across the three scoring variants. Notably, the human-only scoring serves as a ground truth baseline, confirming that the trends observed in automated and calibrated scoring are robust.

These results support the practical use of fully automated scoring for large-scale benchmarking, while human calibration provides additional assurance when reporting final evaluation results.

G.2 Level 1 Analysis

Minor perturbations cause performance drops, revealing shallow generalization. Figure 8 (left) presents model accuracy on Level 1 tasks across four input variants: Original, Perturbed, Knowledge-enhanced, and Math Abstraction. When problems are perturbed through minor changes in wording or numerical values, average model accuracy drops from 82.9% to 81.5%. Notably, Llama 3.3 and Qwen2.5-72B decline by 6.6% and 5.1%, respectively. This indicates that some models exhibit limited robustness and often rely on memorized phrasing or surface patterns rather than generalizable reasoning.

Explicit knowledge prompts mitigate reasoning failures in weaker models. When explicit domain knowledge—such as formulas, constants, or unit conversions—is added to the input, accuracy improves to 85.5% on average. Weaker models benefit the most: GPT-4.1 Mini gains 6.2% and Mixtral-8x7B improves by 9.3%. This pattern suggests that many errors are not caused by a complete lack of knowledge, but rather by the inability to retrieve and apply relevant concepts without targeted prompting. Explicitly embedding domain knowledge thus serves as an effective intervention for enhancing reasoning activation.

Removing contextual language highlights semantic limitations. Performance further increases to 89.4% when problems are rewritten into abstract mathematical form, removing all contextual language. For example, Qwen2.5-7B and Mixtral-8x7B improve by 10.9% and 18.8%, respectively. This reveals that most Level 1 failures are not due to weak computational ability, but rather arise during semantic interpretation and variable binding. Once language ambiguity is removed, models can more reliably execute the required calculations, underscoring a gap between symbolic proficiency and contextual understanding.

Table 2: Level 3 average scores under three scoring variants. Human-calibrated scores are reported as the main results; LLM-only scores are produced by the automated scoring pipeline; Human-only scores are reference scores from official solutions and expert grading.

Model	Human-calibrated		LLM-only		Human-only	
	Original	Perturbed	Original	Perturbed	Original	Perturbed
Human Expert	8.728	8.736	8.697	8.702	8.735	8.729
GPT-4.1	7.108	7.002	7.053	6.972	7.208	7.043
Claude 3.7 Sonnet	6.656	6.445	6.713	6.619	6.970	6.526
GPT-4.1 Mini	6.793	6.378	6.581	6.334	6.705	6.558
DeepSeek-V3	6.289	6.329	6.358	6.264	6.396	6.386
Gemini 2.5 Flash	6.197	6.208	6.002	6.145	6.063	6.185
Gemini 2.0 Flash	6.145	5.991	5.989	5.902	6.167	6.035
GPT-4.1 Nano	5.968	5.738	5.764	5.673	6.074	5.882
GLM-4-32B	5.615	5.653	5.860	5.761	5.760	5.694
GLM-4-9B	4.833	5.169	5.079	5.227	4.822	5.168
Claude 3.5 Sonnet	5.228	5.049	5.317	5.254	5.187	5.106
Llama 3.3	4.837	5.008	4.937	4.804	4.939	5.055
Qwen2.5-72B	4.935	4.722	4.836	4.665	5.007	4.831
Qwen2.5-7B	4.339	4.619	4.580	4.591	4.362	4.669
Llama 4	3.767	3.861	3.808	3.926	3.943	3.892
DeepSeek-R1 7B	4.043	3.810	3.775	3.648	4.105	3.989
Mixtral-8×7B	3.203	3.476	3.110	3.279	3.372	3.577

G.3 Level 2 Analysis

Level 2 tasks emphasize multi-step reasoning under structured constraints, making them more sensitive to input variability. As shown in Figure 8 (right), the average model accuracy declines from 66.6% on the Original version to 61.6% on the Perturbed variant. This 5.0% drop indicates that even minor changes to semantic phrasing or numerical values can significantly disrupt reasoning chains. For instance, GPT-4.1 Nano drops by 9.3% and Qwen2.5-7B by 11.4%, revealing their limited robustness when facing contextual and structural perturbations in problem inputs.

Incorporating explicit domain knowledge helps reduce ambiguity and recover performance. With knowledge-enhanced inputs, the average accuracy rises to 68.6%, a 7.0% improvement over the perturbed baseline. Larger gains are observed for models such as GPT-4.1 Nano (+15.4%) and Qwen2.5-7B (+16.6%), suggesting that knowledge prompts assist in constraint interpretation and formula selection. However, some models such as DeepSeek-V3 show minimal improvement, implying that knowl-

edge access alone may not compensate for limitations in multi-step reasoning capabilities.

Symbolic abstraction of Level 2 tasks into pure mathematical form results in the largest performance gains. The average accuracy increases to 79.2%, with many models gaining over 15%. This trend is especially prominent for weaker models like Qwen2.5-7B (from 37.3% to 69.4%) and Mixtral-8x7B (from 30.0% to 57.7%). These improvements confirm that many model failures stem not from computational weakness, but from difficulties parsing, organizing, and executing the reasoning steps embedded in natural language problem statements. This underscores the importance of assessing upstream cognitive processes that precede symbolic computation—dimensions often underexamined in traditional mathematical benchmarks.

G.4 Level 3 Analysis

Figure 9 presents the performance of various models across four key capabilities: Redundant Information, Multi-Objective Decision, Domain Knowledge, and Uncertainty Handling. The results are

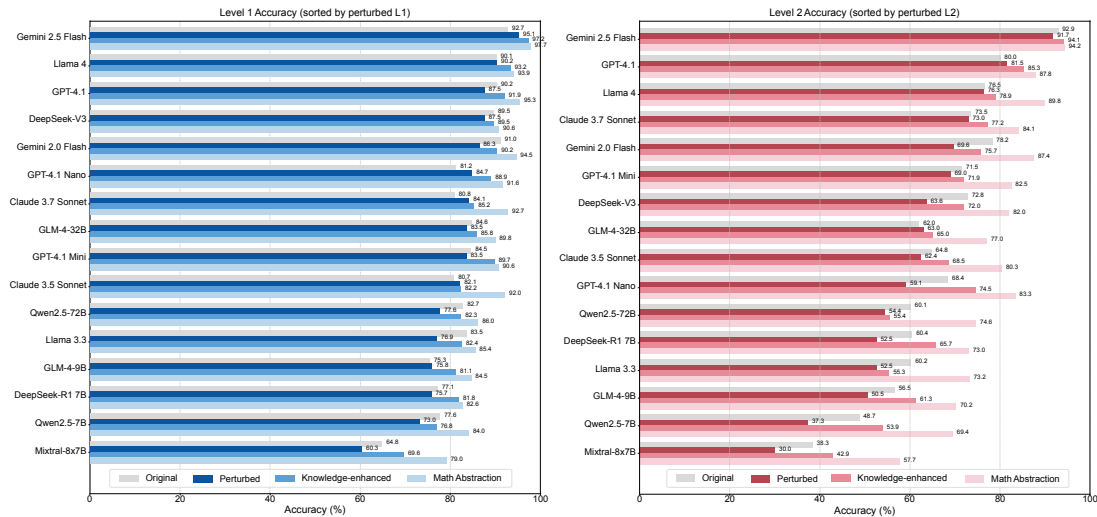


Figure 8: Accuracy of LLMs on Level 1 (left) and Level 2 (right) tasks across four variants: Original, Perturbed, Knowledge-enhanced, and Math Abstraction. Drops in the Perturbed version indicate sensitivity to input changes, while gains in the latter two show that current LLMs require external knowledge or reformulation to improve accuracy—highlighting their lack of these abilities.

further separated into *original* and *perturbed* problem formulations. Overall, human experts substantially outperform all models across all dimensions, with average scores of 8.728 (original) and 8.736 (perturbed). In contrast, LLMs demonstrate significantly lower scores, revealing a persistent gap between current LLMs’ capabilities and human-level reasoning. The average model scores before and after rewriting are 5.372 and 5.341, respectively—a marginal difference of only 0.58%. This indicates that most models possess a reasonable degree of generalization, and the benchmark shows no signs of data contamination across reformulated prompts, preserving task consistency.

Based on the overall average scores, we categorize model performance into three tiers:

Tier 1 (Average Score > 6.5) This tier includes GPT-4.1, Claude 3.7 Sonnet, and GPT-4.1 Mini. These models demonstrate strong performance across all four evaluated capabilities. In particular, their scores in Information Extraction and Multi-Objective Decision often exceed 7, approaching human expert levels. Their performance in Domain Knowledge and Uncertainty Handling also remains consistently above 6, indicating robust reasoning capabilities and broad task adaptability.

Tier 2 (Average Score ≈ 5.5–6.5) This tier consists of DeepSeek-V3, Gemini 2.5 Flash, Gemini 2.0 Flash, GPT-4.1 Nano, and GLM-4-32B. These models achieve reasonable performance in Information Extraction and Multi-Objective Decision, but exhibit noticeable weaknesses in Domain Knowledge and Uncertainty Handling, where scores commonly fall below 6. Some models ap-

proach the 5-point threshold in these dimensions, reflecting limitations in complex reasoning and knowledge integration.

Tier 3 (Average Score < 5.5) This tier includes GLM-4-9B, Claude 3.5 Sonnet, Llama 3.3, Qwen2.5-72B, Qwen2.5-7B, Llama4, DeepSeek-R1 7B, and Mixtral-8x7B. These models consistently underperform across all four capabilities, typically scoring between 3 and 5. Their weakest areas are Domain Knowledge and Uncertainty Handling, where some models fall below 4. These results indicate substantial deficiencies in background reasoning and generalization to ambiguous or under-specified tasks.

G.5 Subfield Performance Analysis

Figures 10 and 11 present an overview of model accuracy across engineering subfields and problem variants for Level 1 and Level 2, respectively.

Model performance varies substantially across engineering subfields. Chemical and biological engineering demonstrates the strongest robustness, with large models maintaining accuracies above 85%, while structural and physical engineering achieves 70–80% and systems and control engineering performs the worst, with large models dropping to 60–70% and small models often below 40%. These results suggest that robustness to contextual perturbations is closely tied to the task characteristics: chemical and biological problems rely more on formulaic knowledge and are less sensitive to input variations, whereas systems and control problems involve more complex reasoning chains and are more vulnerable to perturbations.

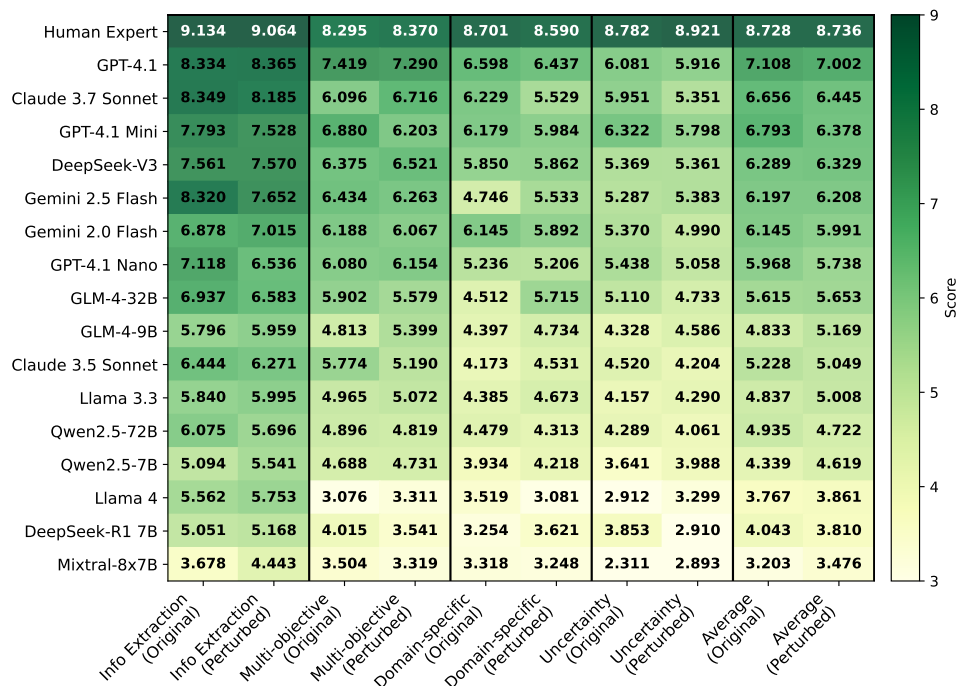


Figure 9: Level 3 Model Evaluation. The figure presents average model performance on Level 3 tasks across four capability dimensions, including information extraction, domain-specific reasoning, multi-objective decision-making, and uncertainty handling, under both original and perturbed problem formulations.

Problem variants reveal subfield-specific differences in knowledge use, reasoning, and robustness, showing that these abilities differ significantly between engineering domains.

The knowledge-enhanced variant substantially improves performance in chemical and biological engineering, moderately benefits structural and physical engineering, and shows limited gains in systems and control engineering, suggesting the latter’s inability to effectively leverage explicit knowledge. Similarly, the math abstraction variant, which isolates mathematical reasoning by removing context, favors chemical and biological engineering, followed by structural and physical engineering, while systems and control engineering remains the weakest. These patterns indicate that the ability to utilize injected knowledge and maintain mathematical reasoning varies considerably across subfields.

structural and physical engineering and negligible gains in systems and control engineering. This indicates that in more complex reasoning and contextual integration tasks, current large language models struggle even more to handle input perturbations, exploit external knowledge effectively, and maintain consistent reasoning, further widening the capability gap across subfields.

The robustness and capability differences across subfields become even more evident under higher task complexity in Level 2.

Compared to Level 1, Level 2 shows larger performance drops under perturbed inputs, highlighting more severe robustness issues. The positive effects of knowledge-enhanced and math abstraction variants remain concentrated in chemical and biological engineering, with only marginal improvements in

2228
2229
2230
2231
2232
2233
2234
2235

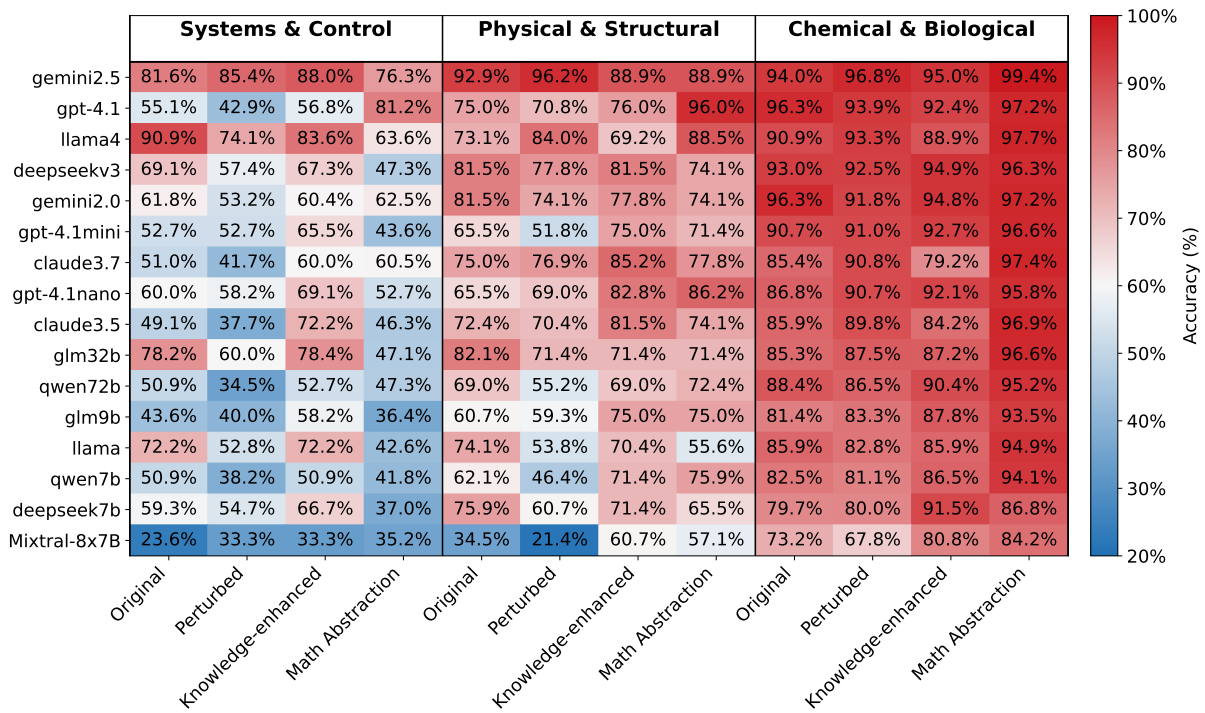


Figure 10: Accuracy across engineering subfields and problem variants in Level 1.

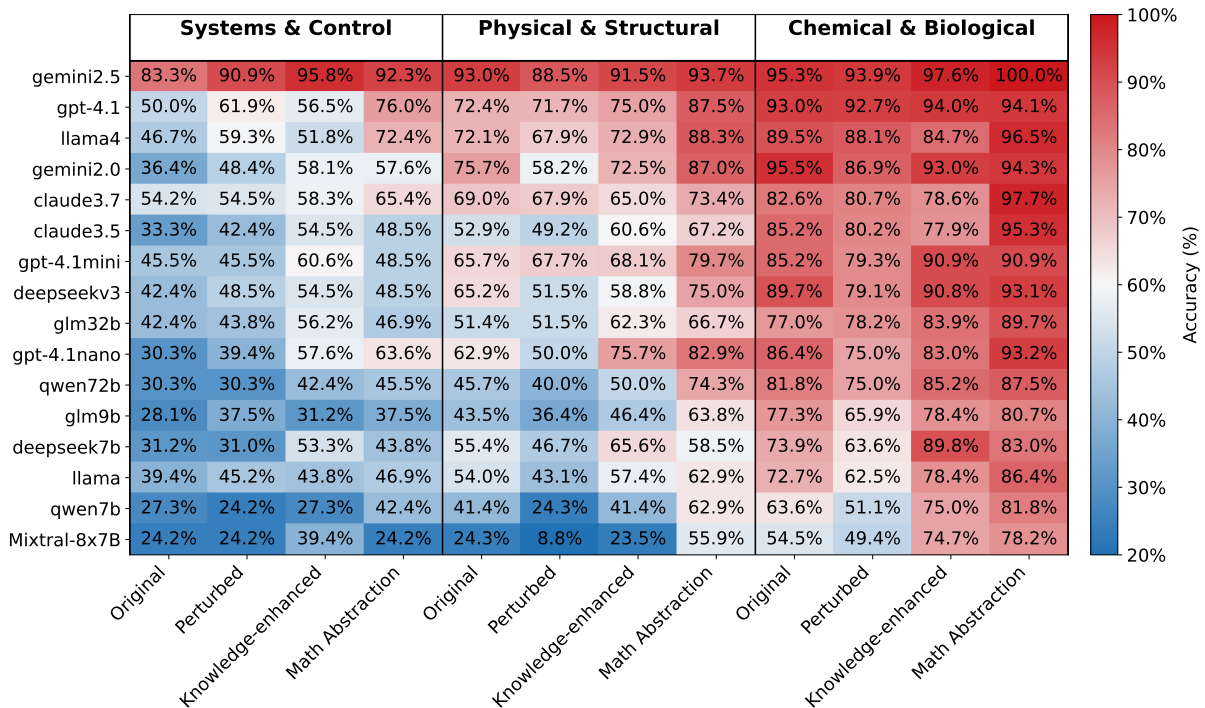


Figure 11: Accuracy across engineering subfields and problem variants in Level 2.