# Solving Inequality Proofs with Large Language Models

**Pan Lu**[*] $^\alpha$, **Jiayi Sheng**[*] $^\beta$, **Luna Lyu**[*] $^\alpha$, **Jikai Jin**$^\alpha$, **Tony Xia**$^\alpha$, **Alex Gu**$^\gamma$, **James Zou**$^\alpha$

$^\alpha$ Stanford University  $^\beta$ UC Berkeley  $^\gamma$ Massachusetts Institute of Technology

🔶 Website: https://ineqmath.github.io/

🐙 Code  🤗 Dataset  🏆 Leaderboard

## Abstract

Inequality proving, crucial across diverse scientific and mathematical fields, tests advanced reasoning skills such as discovering tight bounds and strategic theorem application. This makes it a distinct, demanding frontier for large language models (LLMs), offering insights beyond general mathematical problem-solving. Progress in this area is hampered by existing datasets that are often scarce, synthetic, or rigidly formal. We address this by proposing an *informal yet verifiable* task formulation, recasting inequality proving into two automatically checkable subtasks: bound estimation and relation prediction. Building on this, we release INEQMATH, an expert-curated dataset of Olympiad-level inequalities, including a test set and training corpus enriched with step-wise solutions and theorem annotations. We also develop a novel *LLM-as-judge* evaluation framework, combining a *final-answer* judge with four *step-wise* judges designed to detect common reasoning flaws. A systematic evaluation of 29 leading LLMs on INEQMATH reveals a surprising reality: even top models like o1 achieve less than 10% overall accuracy under step-wise scrutiny; this is a drop of up to 65.5% from their accuracy considering only final answer equivalence. This discrepancy exposes fragile deductive chains and a critical gap for current LLMs between merely finding an answer and constructing a rigorous proof. Scaling model size and increasing test-time computation yield limited gains in overall proof correctness. Instead, our findings highlight promising research directions such as theorem-guided reasoning and self-refinement.
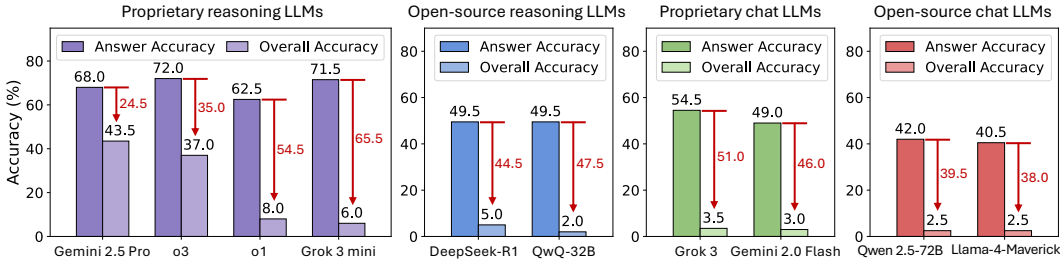
Figure 1: Final-answer accuracy versus overall accuracy for leading LLMs across different categories on the INEQMATH benchmark of Olympiad-level inequality problems. Overall accuracy, measuring both answer correctness and step soundness, is substantially lower than final-answer accuracy for all model types. This highlights a critical gap: while LLMs may find correct final answers to these inequality problems, their reasoning is often unsound. Each model used its optimal maximal tokens.

---

* Co-first authors. Corresponding authors: ✉{panlu, jamesz}@stanford.edu

# 1 Introduction

Mathematical inequalities are fundamental to diverse fields such as analysis, optimization, and probability theory, with applications spanning scientific modeling, economics, and competitive mathematics. Proving an inequality is a complex endeavor, demanding not just calculation but a sophisticated blend of intuition for discovering tight bounds, strategic insight for selecting and applying classical theorems (e.g., AM-GM, Cauchy-Schwarz), and precise symbolic transformations. These skills are hallmarks of advanced mathematical reasoning, distinguishing inequality proving from general math problem-solving. Automating this process would therefore have broad impact: it could supply automated theorem provers (ATPs) with missing lemmas, accelerate formal verification processes, and serve as a demanding testbed for general-purpose reasoners. However, despite impressive advancements in LLMs like DeepSeek-R1 [14] and OpenAI o3 [48], as well as in ATPs themselves [16, 18, 26, 34, 50, 75], automating inequality proving remains a challenging frontier.

A major bottleneck in advancing LLM capabilities for inequality proving is the scarcity of suitable benchmarks. Existing resources fall short in several ways: general ATP collections like MiniF2F [82] and ProofNet [7] contain few inequalities; synthetic datasets such as INT [64] and AIPS [63] offer scale but may lack structural diversity due to template-based generation; and curated collections like ChenNEQ [8] are often too small for extensive training. More fundamentally, most existing datasets adopt a *fully formal* representation, where problems and proofs are encoded in systems like Lean [11] or Isabelle [42]. While formal mathematical reasoning offers correctness guarantees and is a vital research direction, LLMs, trained on vast corpora of natural language, often exhibit strong informal reasoning capabilities. This suggests LLMs might solve problems informally even when struggling with the exacting syntax of formal provers. Our work, therefore, aims to explore and benchmark these informal abilities, complementing formal mathematical AI by focusing on a mode of reasoning closer to human intuition and the preliminary, often less structured, stages of mathematical discovery.

To bridge this gap between formal rigor and intuitive problem-solving, we propose an *informal yet verifiable* formulation (§2). Rather than requiring fully machine-checkable proofs within formal systems, we reformulate inequality problems into two concrete, automatically verifiable subtasks: (i) *Bound estimation*—determine the largest (or smallest) constant $C$ that preserves the inequality; and (ii) *Relation prediction*—identify which relation ($>, \geq, =, \leq,$ or $<$) holds between two expressions. Both tasks can be presented in natural language and LaTeX, solved step-by-step by an LLM, and their final answers (a constant or a relation symbol) can be automatically checked. This preserves the creative essence of inequality proving while avoiding the heavy overhead of formal proof assistants.

Building on this formulation, we present INEQMATH (§3), the first large-scale dataset of Olympiad-level inequalities written entirely in informal language. The *test set* comprises 200 original problems, each crafted and reviewed by IMO-level medalists to ensure both originality and difficulty. The *training corpus* includes 1,252 problems sourced from advanced textbooks, automatically rephrased by LLMs into our subtasks and then meticulously reviewed by human experts. A key feature is that each training problem is accompanied by up to four *step-wise solution paths*, providing rich data for training LLMs on fine-grained reasoning. Additionally, 76.8% of the training problems are annotated with 83 *named theorems* across 29 categories relevant to their solutions. As shown in Table 2, INEQMATH surpasses prior resources in scale, diversity, and alignment with human-like, informal problem-solving approaches.

However, producing the correct final answer is insufficient; the reasoning process itself must be sound. To rigorously assess this, we introduce an *LLM-as-judge* evaluation framework (§4). This framework comprises a high-precision *final-answer judge* to verify the answer equivalence, complemented by four specialized *step-wise judges* for step soundness. These step-wise judges are designed to detect the frequent reasoning flaws identified in our pilot studies: inappropriate reliance on *toy case* examples, unaddressed *logical gaps*, unjustified *numeric approximations*, and *numeric calculation* errors. Validated on manually labeled development set solutions, these judges demonstrate high reliability (F1 > 0.9 on average) and offer a scalable method to scrutinize the deductive integrity of LLM-generated proofs.

We evaluate 29 leading LLMs ranging from chat models to advanced reasoning LLMs, both open-source and proprietary (§5). As key results highlighted in Figure 1, several key findings emerge. While specialized reasoning LLMs (e.g., o1 [45]) achieve higher *final-answer* accuracy than general-purpose chat models (e.g., GPT-4o [43]), this advantage often collapses under step-wise scrutiny. Once our judges inspect every reasoning step, *overall* accuracy plummets by up to 65.5%. Indeed, even

top-performing models like o1 achieve less than 10% overall accuracy (Table 4), exposing fragile deductive chains and a significant gap between finding an answer and constructing a rigorous proof.

Our in-depth study (§5.3) reveals that while larger model sizes correlate with improved *final-answer accuracy*, their impact on *overall accuracy* is limited (e.g., o1 achieves only 8.0% overall accuracy). Similarly, extending test-time computation through longer reasoning chains offers diminishing returns in overall correctness (e.g., o1's 8.0% overall accuracy remains unchanged when scaling maximum completion tokens from 5K to 40K, while o3 [48] saturates around 31%). These findings suggest that current scaling approaches are insufficient for robust deductive reasoning in IneqMath. Instead, we explore promising improvement strategies, demonstrating potential gains from methods such as theorem-guided reasoning—by providing golden theorems (improving overall accuracy by up to 11% for o3-mini [47] ) and critic-guided self-refinement (e.g., a 5% absolute increase in overall accuracy for Gemini 2.5 Pro [22]).

In summary, our work makes four key contributions: 1) We introduce an *informal* reformulation of inequality proving, decomposing the task into two verifiable subtasks (§2). 2) We release IneqMath, an expert-curated benchmark of Olympiad-level inequalities and a training corpus enriched with step-wise solutions and theorem annotations (§3). 3) We develop a modular *LLM-as-judge* framework that rigorously evaluates both final answers and proof step soundness (§4). 4) We conduct a systematic empirical study (§5) that exposes a pronounced gap between LLM performance and mathematical rigor, highlighting avenues for future research.

## 2 Task Formalization: An Informal Perspective

Inequality proof problems require demonstrating that a specified inequality holds under given conditions, such as proving $a + b \geq 2\sqrt{ab}$ for all positive real numbers $a$ and $b$. Traditionally, these problems are formalized in proof assistants like Lean or Isabelle, represented as a tuple $(S_0, I, P)$, where $S_0$ is the initial state, $I$ is the inequality, and $P$ is a set of premises. The proof process, often modeled as a Markov Decision Process, constructs a step-by-step solution verified by the system. However, this formal approach demands expertise in specialized tools, while informal proofs in natural language, though more intuitive, are difficult to verify automatically due to their unstructured nature.

To address these challenges, we propose an *informal* perspective that reformulates inequality proof problems into two *verifiable* subtasks: **bound estimation** and **relation prediction**.

---

**IneqMath Training Example 1: Bound Problem**

**Question:** Find the maximal constant $C$ such that for all real numbers $a, b, c$, the inequality holds:

$$\sqrt{a^2 + (1-b)^2} + \sqrt{b^2 + (1-c)^2} + \sqrt{c^2 + (1-a)^2} \geq C$$

**Solution:** Applying Minkowsky's Inequality to the left-hand side we have

$$\sqrt{a^2 + (1-b)^2} + \sqrt{b^2 + (1-c)^2} + \sqrt{c^2 + (1-a)^2} \geq \sqrt{(a+b+c)^2 + (3-a-b-c)^2}$$

By denoting $a + b + c = x$, we get

$$\sqrt{(a+b+c)^2 + (3-a-b-c)^2} = \sqrt{2\left(x - \frac{3}{2}\right)^2 + \frac{9}{2}} \geq \sqrt{\frac{9}{2}} = \boxed{\frac{3\sqrt{2}}{2}}.$$

**Minkowsky's Inequality Theorem:** For any real number $r \geq 1$ and any positive real numbers $a_1, a_2, \ldots, a_n, b_1, b_2, \ldots, b_n$

$$\left(\sum_{i=1}^{n} (a_i + b_i)^r\right)^{\frac{1}{r}} \leq \left(\sum_{i=1}^{n} a_i^r\right)^{\frac{1}{r}} + \left(\sum_{i=1}^{n} b_i^r\right)^{\frac{1}{r}}$$

---

This **bound estimation** task involves finding an optimal constant for a given inequality. For example, in $a + b \geq C\sqrt{ab}$ for $\forall a, b > 0$, the objective is to find the largest $C$. Formally, a bound estimation problem instance is a triple:

$$\Pi_{\text{bound}} = \left(f(\mathbf{x}),\ g(\mathbf{x}),\ \mathcal{D}\right), \quad \text{where } \mathcal{D} \subseteq \mathbb{R}^n.$$

Here, $f, g : \mathcal{D} \to \mathbb{R}$ are two expressions involving variables $\mathbf{x} = (x_1, \ldots, x_n)$ within a specified domain $\mathcal{D}$ (e.g., $x_i > 0$, $\sum x_i = 1$), and $g(\mathbf{x}) > 0, \forall \mathbf{x} \in \mathcal{D}$. The goal is to determine the extremal:

$$C^\star = \sup\{C \in \mathbb{R} : f(\mathbf{x}) \geq Cg(\mathbf{x}), \forall \mathbf{x} \in \mathcal{D}\} \text{ or } C^\star = \inf\{C \in \mathbb{R} : f(\mathbf{x}) \leq Cg(\mathbf{x}), \forall \mathbf{x} \in \mathcal{D}\}.$$

The **relation prediction** task requires determining the correct relationship between two expressions. For instance, given expressions $f(\mathbf{x}) = a + b$ and $g(\mathbf{x}) = 2\sqrt{ab}$, the goal is to identify the relation (in this case, $\geq$) that holds for $\forall a, b > 0$. Formally, a relation prediction problem instance is a triple:

$$\Pi_{\text{rel}} = \big( f(\mathbf{x}), \, g(\mathbf{x}), \, \mathcal{D} \big),$$

where $f(\mathbf{x})$ and $g(\mathbf{x})$ are expressions over variables $\mathbf{x}$ in domain $\mathcal{D} \subseteq \mathbb{R}^n$. The goal is to find the relation between $f(\mathbf{x})$ and $g(\mathbf{x})$ (i.e. $>, \geq, =, \leq, <$, or none of the above).

---

**INEQMATH Training Example 2: Relation Problem**

**Question:** Let $a, b, c$ be positive real numbers such that $abc = 1$. Consider the following expressions:

$$\frac{b+c}{\sqrt{a}} + \frac{c+a}{\sqrt{b}} + \frac{a+b}{\sqrt{c}} \quad (\quad) \quad \sqrt{a} + \sqrt{b} + \sqrt{c} + 3$$

Determine the correct inequality relation to fill in the blank.

Options: (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

**Solution:** From the AM-GM Inequality, we have

$$\frac{b+c}{\sqrt{a}} + \frac{c+a}{\sqrt{b}} + \frac{a+b}{\sqrt{c}} \geq 2 \left( \sqrt{\frac{bc}{a}} + \sqrt{\frac{ca}{b}} + \sqrt{\frac{ab}{c}} \right)$$

$$= \left( \sqrt{\frac{bc}{a}} + \sqrt{\frac{ca}{b}} \right) + \left( \sqrt{\frac{ca}{b}} + \sqrt{\frac{ab}{c}} \right) + \left( \sqrt{\frac{ab}{c}} + \sqrt{\frac{bc}{a}} \right)$$

$$\geq 2(\sqrt{a} + \sqrt{b} + \sqrt{c}) \boxed{\geq} \sqrt{a} + \sqrt{b} + \sqrt{c} + 3\sqrt[6]{abc} = \sqrt{a} + \sqrt{b} + \sqrt{c} + 3.$$

**AM-GM Inequality Theorem:** If $a_1, a_2, \ldots, a_n$ are nonnegative real numbers, then

$$\frac{1}{n} \sum_{i=1}^{n} a_i \geq (a_1 a_2 \ldots a_n)^{\frac{1}{n}}$$

with equality if and only if $a_1 = a_2 = \ldots = a_n$. This is a special case of the Power Mean Inequality.

---

These subtasks are chosen because they frequently appear in mathematical problem-solving, simplify the evaluation process, and crucially, retain the core reasoning challenges inherent in original inequality proof problems. An ideal LLM solution should not only produce the correct final answer but also present a clear, logically sound, and complete derivation. This includes strategic application of theorems, accurate symbolic manipulations and calculations, and justification of all critical steps.

## 3 INEQMATH: The Inequality Problem Dataset

This section describes the data curation process and key statistics of INEQMATH, a novel collection of inequality problems designed to support the informal perspective on solving and proving inequalities.

**Test data curation.** To mitigate contamination from common sources (textbooks, contests and online resources) that may be present in LLM training corpora, we commissioned IMO-level medalists to design novel inequality problems. These underwent rigorous review by a separate expert group and were validated only upon unanimous confirmation of solvability, soundness, and ground truth correctness. Problems identified as easier by experts were excluded from the test set (repurposed for development) to ensure a high level of challenge. To further illustrate the modest contamination, we also conduct a memorization probe on the test set in §C.8. See the developed curation tool in §A.2. We host an online evaluation website[1], providing a fair evaluation platform for the community.

**Key statistics.** As shown in Table 1, the INEQMATH dataset comprises 200 test problems for benchmarking, 100 development problems with public ground truth, and 1,252 training problems split evenly between bound estimation and relation prediction tasks. Each training problem includes step-wise solutions, with up to four solutions per problem, and 76.8% (962 problems) are annotated with relevant theorems. The dataset features 83 named theorems across 29 categories, with their distribution illustrated in Figure 2. Test problem examples are provided in §A.4.

---

[1] https://huggingface.co/spaces/AI4Math/IneqMath-Leaderboard

| Statistic | Number | Bnd. | Rel. |
|---|---|---|---|
| Theorem categories | 29 | - | - |
| Named theorems | 83 | - | - |
| Training problems (for training) | 1252 | 626 | 626 |
| - With theorem annotations | 962 | 482 | 480 |
| - With solution annotations | 1252 | 626 | 626 |
| - Avg. solutions per problem | 1.05 | 1.06 | 1.05 |
| - Max solutions per problem | 4 | 4 | 4 |
| Dev problems (for development) | 100 | 50 | 50 |
| Test problems (for benchmarking) | 200 | 96 | 104 |

Table 1: Statistics of the INEQMATH dataset.



Figure 2: Distribution of theorem categories.

**Comparison to existing datasets.** As summarized in Table 2, INEQMATH stands out for: (1) providing expert-curated training and test sets, (2) offering rich annotations with step-wise solutions and 83 grounded theorems, and (3) adopting an informal, accessible format for inequality proving through bound estimation and relation prediction, evaluated via *LLM-as-judge*. This design bridges the gap between formal proof systems and intuitive mathematical reasoning, making INEQMATH a unique resource for advancing LLM capabilities in problem solving and theorem proving.

| Datasets | Data Source | | Data Annotation | | Problem and Evaluation | | |
|---|---|---|---|---|---|---|---|
| | **Training** | **Test / Dev** | **#Theorem** | **Solution** | **Category** | **Format** | **Evaluation** |
| INT [64] | Synthesized | Synthesized | 35 | ✓ | Proof | Formal | Symbolic DSL |
| AIPS [63] | Synthesized | ✗ | 8 | ✓ | Proof | Formal | Symbolic DSL |
| MO-INT [63] | ✗ | Data compilation | ✗ | ✗ | Proof | Formal | Symbolic DSL |
| MINIF2F [82] | ✗ | Autoformalization | ✗ | ✗ | Proof | Formal | L∃∀N |
| ProofNet [7] | ✗ | Autoformalization | ✗ | ✗ | Proof | Formal | L∃∀N |
| FormalMATH [77] | ✗ | Autoformalization | ✗ | ✗ | Proof | Formal | L∃∀N |
| leanWorkbook [76] | Autoformalization | Autoformalization | ✗ | ✗ | Proof | Formal | L∃∀N |
| Proof or Bluff [49] | ✗ | Data compilation | ✗ | ✗ | Proof | Informal | Human judge |
| CHAMP [39] | ✗ | Autoformalization | ✗ | ✗ | Open | Informal | Human judge |
| Putnam Axiom [23] | ✗ | Data compilation | ✗ | ✗ | Open | Informal | Answer checking |
| LiveMathBench [37] | ✗ | Data compilation | ✗ | ✗ | Open | Informal | Answer checking |
| **INEQMATH (Ours)** | Expert annotated | Expert annotated | 83 | ✓ | MC, Open | Informal | LLM-as-judge |

Table 2: Comparison of datasets for inequality and theorem proving. INEQMATH provides expert-annotated training and test/dev sets, featuring high-quality named theorems and step-wise solutions for model development. Unlike prior datasets using synthesis or autoformalization, INEQMATH presents problems in informal language across multiple-choice (MC) and open-ended (Open) formats, and employs *LLM-as-judge* for evaluation.

**Potential contamination statement.** To ensure rigorous evaluation, the INEQMATH test set was commissioned from IMO-level medalists to feature novel problems, minimizing prior LLM pre-training exposure. The poor performance across models (§5.2), particularly in overall accuracy (which demands step-wise correctness), strongly suggests that the benchmark poses a significant reasoning challenge, regardless of any potential familiarity with the underlying mathematical concepts. We therefore believe the INEQMATH test set effectively probes novel problem-solving capabilities, and our conclusions on current LLM limitations in rigorous inequality proving remain robust.

## 4 Fine-grained Informal Judges for Inequality Solving

The test split of the INEQMATH dataset serves as our benchmark, comprising 200 Olympiad-level inequality problems that challenge both humans and current LLMs. Traditional evaluation methods fall short in this setting: expert annotation is accurate but prohibitively labor-intensive, while automated techniques such as string matching or value equivalence fail to capture step-by-step correctness—an essential aspect of inequality problem solving. To address this, we propose a fine-grained *LLM-as-judge* framework as illustrated in Figure 3, consisting of a *final-answer judge* for verifying the predicted answer (§4.1) and four specialized *step-wise judges* targeting common reasoning flaws (§4.2). A solution is considered correct *overall* only if it passes all five judges. As shown in Table 3, these judges achieve strong alignment with human annotations (F1 = 0.93), providing a scalable yet reliable alternative to manual evaluation.

### 4.1 Final Answer Judge

LLM-generated solutions to INEQMATH problems typically involve multiple reasoning steps followed by a concluding answer statement. However, the final answer may vary in phrasing, structure, or numeric format, especially for bound estimation problems. For example, $C = \frac{1}{\sqrt{2}}$ and $C = \frac{\sqrt{2}}{2}$
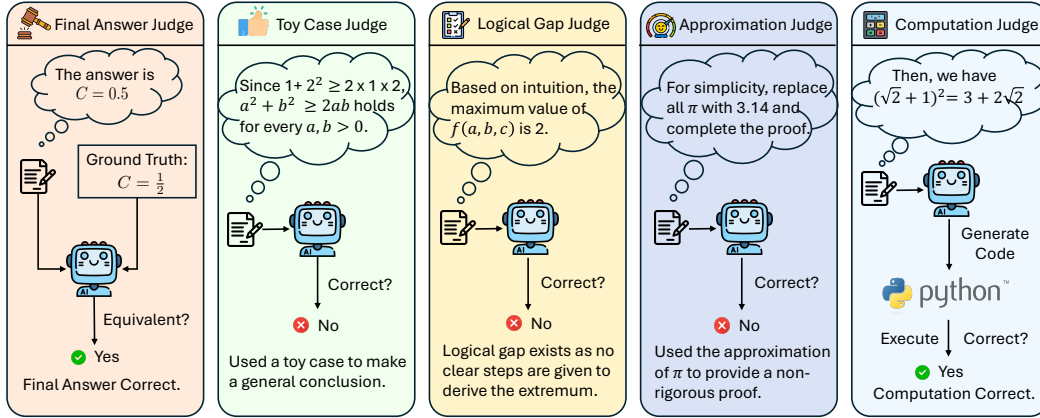
5

Figure 3: Illustration of the fine-grained *LLM-as-judge* framework. The framework combines a *Final Answer Judge* with four step-wise judges: *Toy Case Judge*, *Logical Gap Judge*, *Numerical Approximation Judge* (shown as *Approximation Judge*), and *Numerical Computation Judge* (shown as *Computation Judge*). A solution is considered correct only if it passes all five judges.

| LLM-as-Judge | Judge type | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|---|
| Final Answer Judge | Answer checking | 1.00 | 1.00 | 1.00 | 1.00 |
| Toy Case Judge | Step soundness | 0.91 | 0.86 | 0.97 | 0.91 |
| Logical Gap Judge | Step soundness | 0.96 | 0.95 | 0.98 | 0.96 |
| Numerical Approximation Judge | Step soundness | 0.96 | 0.95 | 0.98 | 0.96 |
| Numerical Computation Judge | Step soundness | 0.71 | 0.68 | 0.98 | 0.80 |
| **Average** | - | 0.91 | 0.89 | 0.98 | **0.93** |

Table 3: Performance metrics of *LLM-as-judge* framework on development set.

are mathematically equivalent but differ in form. Recent work [38] evaluates LLM outputs via format normalization and exact string matching, without accounting for mathematical equivalence. To address this, we propose a two-stage *Final Answer Judge*: it first identifies the concluding sentence containing the predicted answer, and then performs robust equivalence checking to assess mathematical correctness, even when the form differs from the reference. Prompt details and examples are in §B.1.

## 4.2 Four Step-wised Judges

**Toy Case Judge.** Inequality problems in INEQMATH often require reasoning over continuous domains (e.g., all $a, b, c > 0$), where specific numerical examples alone are insufficient for a valid proof. LLM frequently generalizes incorrectly from such examples—e.g., claiming an inequality holds universally because it holds for $a = 1, b = 2, c = 3$. Prior work [17] flags these under a broad "logical flaw" category, lacking granularity for targeted analysis. Our *Toy Case Judge* addresses this by detecting unjustified generalization from toy examples. It prompts an LLM to flag conclusions based solely on specific instances without broader justification. See §B.2 for prompts and examples.

**Logical Gap Judge.** INEQMATH inequality problems often involve multi-step derivations (e.g., algebraic manipulation, constrained optimization, functional transformations) needing explicit justification. LLMs, however, often skip key reasoning steps or assert conclusions without support (e.g., stating an optimal bound without derivation). Existing step-level evaluations [68] assess validity and redundancy but lack granularity for such logical omissions. Our *Logical Gap Judge* addresses this by flagging missing transitions, unjustified claims, and vague derivations, especially in steps involving inequality transformations or bound estimation (see §B.3 for details).

**Numerical Approximation Judge.** Inequality problems in INEQMATH often demand exact symbolic reasoning, where the use of numeric approximations—e.g., replacing $\sqrt{2}$ with $1.414$—can compromise mathematical rigor. However, many LLM-generated solutions resort to such approximations in intermediate steps, leading to inaccurate or non-generalizable conclusions. To address this, we introduce a *Numerical Approximation Judge* that flags inappropriate use of numeric approximations—specifically when they affect derivations or final answers. Approximations used solely for intuition or side remarks are permitted. See §B.4 for prompt details and examples.
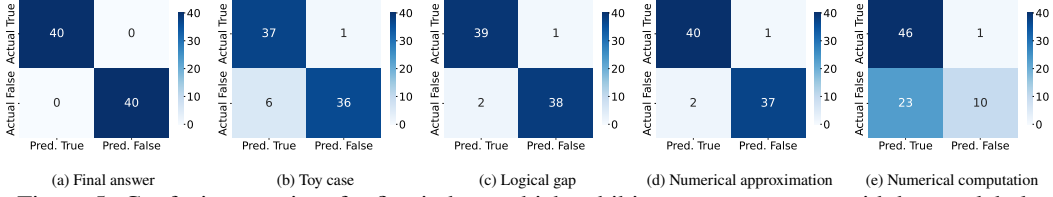
6

Figure 5: Confusion matrices for five judges, which exhibit strong agreement with human labels.

**Numerical Computation Judge.** Many INEQMATH problems require explicit numerical computations after variable assignment (e.g., evaluating $\frac{27}{2}$ or summing rational terms). While symbolic reasoning is vital, arithmetic accuracy is equally crucial for overall correctness. Prior work (e.g., EIC-Math [30]) categorizes broad error types but often overlooks subtle miscalculations in multi-step derivations. Our *Numerical Computation Judge* addresses this by verifying arithmetic steps once variables are instantiated. It prompts an LLM to extract numerical expressions, convert them into Python code, and evaluate using floating-point arithmetic within a small tolerance. This enables high-precision checking of both intermediate and final results. To further improve precision and mitigate floating-point issues, we encourage the use of symbolic mathematics packages such as SymPy, particularly for handling fractions and decimal numbers. Additional details are provided in §B.5.

### 4.3 Effectiveness Verification of Judges

**A holistic LLM judge baseline.** To motivate our specialized judging system, we first evaluate a heuristic *LLM-as-judge* baseline. This prompts a single, general-purpose LLM to holistically assess INEQMATH solution correctness, based on both final answer accuracy and step-wise soundness across the four reasoning categories in §4.2. As shown in the confusion matrix (Figure 4) using 80 human-annotated development examples, this naive approach exhibits poor agreement with human labels, underscoring its unreliability for rigorous evaluation in this domain.
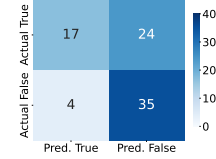


Figure 4: Confusion matrix for the judge baseline.

**Performance of our fine-grained judges.** In contrast, our proposed suite of five specialized judges exhibits strong alignment with human evaluations. Figure 5 presents the confusion matrices for each judge on the same development set. The final answer judge (using GPT-4o-mini) achieves near-perfect agreement, while the four step-wise judges (chosen for a balance of performance and cost as detailed in §B.6) also demonstrate high fidelity. This confirms that decomposing the complex evaluation task into targeted sub-problems allows LLMs to serve as reliable evaluators.

**Qualitative analysis of judge disagreements.** Despite the strong aggregate performance (overall F1 = 0.93, Table 3), *LLM-as-judge* evaluations are not perfect. Acknowledging the skepticism surrounding LLM-based evaluation, we conducted a qualitative analysis of failure cases where our judges' assessments diverged from human annotations. Detailed examples are provided in §B.7. These instances underscore that while highly effective, our LLM judges can still struggle with the deep, nuanced understanding that characterizes expert-level mathematical reasoning.

## 5 Experiments in INEQMATH

### 5.1 Experimental Setups

We conduct a systematic evaluation of 29 leading LLMs on the inequality problems in the INEQMATH test set. The evaluated models span two categories: general-purpose chat models (both open-source and proprietary) and specialized reasoning LLMs designed for complex, multi-step problem-solving. All models are prompted in a zero-shot setting with the problem statement and the instruction: "*Please solve the problem with clear, rigorous, and logically sound steps*" to encourage detailed reasoning. Model responses are assessed using our *LLM-as-judge* framework (§4). We report three key metrics:

- *Answer Acc*: Measures the predicted answer correctness, verified by the final-answer judge (§4.1).
- *Step Acc*: Aggregates the correctness of individual reasoning steps as determined by our four specialized step-wise judges (§4.2), which target common flaws.
- *Overall Acc*: The primary metric, which deems a solution correct only if it achieves both a correct final answer and flawless step-wise reasoning (i.e., passes all five judges).

A response is thus considered fully correct (*Overall Acc*) only if it produces a correct final answer through logically valid steps, passing scrutiny from all judges. Additional setup details are in §C.1.

| Models | Overall Acc (↑) | | | Answer Acc (↑) | | | Step Acc (↑) | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | No Toy Case | | | No Logic. Gap | | | No Approx. Error | | | No Comp. Error | | | | | |
| | All | Bnd. | Rel. | All | Bnd. | Rel. | All | Bnd. | Rel. | All | Bnd. | Rel. | All | Bnd. | Rel. | All | Bnd. | Rel. | | | |
| *Heuristic Methods* | | | | | | | | | | | | | | | | | | | | | |
| Random Guess | - | - | - | 8.5 | 0.0 | 16.3 | - | - | - | - | - | - | - | - | - | - | - | - | | | |
| Frequent Guess | - | - | - | 18.0 | 9.4 | 26.0 | - | - | - | - | - | - | - | - | - | - | - | - | | | |
| *Open-source Chat LLMs* | | | | | | | | | | | | | | | | | | | | | |
| Qwen2.5-Coder-32B [27] | $1.5_{39.0}$ | $1.0_{50.0}$ | $1.9_{28.9}$ | 40.5 | $\underline{51.0}$ | 30.8 | 36.0 | 27.1 | 44.2 | 3.0 | 2.1 | 3.8 | 90.5 | $\underline{96.9}$ | 84.6 | 88.5 | 89.6 | 87.5 | | | |
| Llama-4-Scout [41] | $1.5_{32.0}$ | $2.1_{44.8}$ | $1.0_{20.2}$ | 33.5 | 46.9 | 21.2 | 30.5 | 15.6 | 44.2 | 3.5 | $\underline{4.2}$ | 2.9 | $\underline{93.0}$ | 94.8 | $\underline{91.3}$ | 92.5 | 92.7 | 92.3 | | | |
| Qwen2.5-72B [51] | $2.5_{39.5}$ | $3.1_{47.9}$ | $1.9_{31.8}$ | 42.0 | $\underline{51.0}$ | 33.7 | $\underline{54.5}$ | $\underline{53.1}$ | $\underline{55.8}$ | 5.0 | $\underline{4.2}$ | $\underline{5.8}$ | 91.0 | 94.8 | 87.5 | $\underline{95.0}$ | $\underline{94.8}$ | 95.2 | | | |
| Llama-4-Maverick [40] | $2.5_{38.0}$ | $2.1_{43.7}$ | $2.9_{32.7}$ | 40.5 | 45.8 | $\underline{35.6}$ | 42.5 | 28.1 | $\underline{55.8}$ | 4.0 | $\underline{4.2}$ | 3.8 | 89.0 | 91.7 | 86.5 | $\underline{95.0}$ | 92.7 | $\underline{97.1}$ | | | |
| Qwen2.5-7B [52] | $3.0_{32.0}$ | $2.1_{38.5}$ | $3.8_{26.0}$ | 35.0 | 40.6 | 29.8 | 44.5 | 32.3 | $\underline{55.8}$ | 4.5 | 3.1 | $\underline{5.8}$ | 92.5 | $\underline{96.9}$ | 88.5 | 93.0 | 92.7 | 93.3 | | | |
| *Proprietary Chat LLMs* | | | | | | | | | | | | | | | | | | | | | |
| Gemini 2.0 Flash-Lite [20] | $1.5_{31.5}$ | $2.1_{41.7}$ | $1.0_{22.1}$ | 33.0 | 43.8 | 23.1 | 11.5 | 11.5 | 11.5 | 3.5 | 3.1 | 3.8 | 73.0 | 77.1 | 69.2 | 90.5 | 87.5 | 93.3 | | | |
| GPT-4o mini [44] | $2.0_{37.5}$ | $1.0_{41.7}$ | $2.9_{33.6}$ | 39.5 | 42.7 | 36.5 | 29.0 | 11.5 | $\underline{45.2}$ | 2.5 | 2.1 | 2.9 | 90.0 | 91.7 | 88.5 | 93.0 | 92.7 | 93.3 | | | |
| GPT-4.1 [46] | $2.5_{38.0}$ | $0.0_{31.3}$ | $4.8_{44.2}$ | 40.5 | 31.3 | $\underline{49.0}$ | 16.0 | 12.0 | 19.0 | 10.0 | 8.3 | 11.5 | 59.5 | 66.7 | 52.9 | 93.5 | 92.7 | $\underline{94.2}$ | | | |
| GPT-4o [43] | $3.0_{34.5}$ | $2.1_{38.5}$ | $3.8_{30.8}$ | 37.5 | 40.6 | 34.6 | $\underline{32.0}$ | $\underline{21.9}$ | 43.0 | 3.5 | 3.1 | 3.8 | $\underline{92.5}$ | $\underline{93.8}$ | $\underline{91.4}$ | 94.0 | 93.8 | $\underline{94.2}$ | | | |
| Gemini 2.0 Flash [19] | $3.0_{46.0}$ | $3.1_{56.3}$ | $2.9_{36.5}$ | 49.0 | 59.4 | 39.4 | 15.5 | 13.5 | 17.3 | 13.5 | 7.3 | 19.2 | 55.5 | 60.4 | 51.0 | $\underline{94.5}$ | 94.8 | $\underline{94.2}$ | | | |
| Grok 3 [66] | $3.5_{51.0}$ | $4.2_{62.5}$ | $2.9_{40.4}$ | $\underline{54.5}$ | $\underline{66.7}$ | 43.3 | 17.0 | 13.7 | 20.2 | $\underline{16.0}$ | $\underline{11.6}$ | $\underline{20.2}$ | 36.0 | 42.1 | 30.8 | 93.0 | $\underline{96.8}$ | 90.4 | | | |
| *Open-source Reasoning LLMs* | | | | | | | | | | | | | | | | | | | | | |
| QwQ-32B [5] | $2.0_{47.5}$ | $2.1_{52.1}$ | $1.9_{43.3}$ | 49.5 | 54.2 | 45.2 | 26.0 | 25.0 | 26.9 | 29.5 | 20.1 | 37.5 | 21.0 | 20.8 | 21.2 | 87.0 | 82.3 | $\underline{91.3}$ | | | |
| Deepseek-R1 (Llama-70B) [12] | $3.5_{50.0}$ | $5.2_{53.1}$ | $1.9_{47.1}$ | $\underline{53.5}$ | $\underline{58.3}$ | $\underline{49.0}$ | 23.0 | 24.0 | 22.1 | 26.0 | 20.9 | 30.8 | $\underline{35.5}$ | $\underline{38.5}$ | 32.7 | 87.0 | 89.6 | 84.6 | | | |
| Deepseek-R1 (Qwen-14B) [13] | $5.0_{35.5}$ | $6.3_{36.4}$ | $3.8_{34.7}$ | 40.5 | 42.7 | 38.5 | 21.0 | 18.8 | 23.1 | 21.0 | 19.8 | 22.1 | $\underline{35.5}$ | $\underline{38.5}$ | 32.7 | 85.0 | 91.7 | 78.8 | | | |
| Deepseek-R1 [14] | $5.0_{44.5}$ | $4.2_{63.5}$ | $5.8_{26.9}$ | 49.5 | 67.7 | 32.7 | 57.0 | 53.1 | 60.9 | 17.5 | 6.3 | 27.9 | 81.0 | 95.8 | 67.3 | 95.0 | $\underline{99.0}$ | $\underline{91.3}$ | | | |
| Qwen3-235B-A22B [53] | $6.0_{35.0}$ | $3.1_{32.3}$ | $8.7_{37.5}$ | 41.0 | 35.4 | 46.2 | $\underline{35.0}$ | $\underline{30.2}$ | $\underline{39.4}$ | $\underline{36.0}$ | $\underline{26.0}$ | $\underline{45.2}$ | 31.0 | 28.1 | $\underline{33.7}$ | 92.5 | $\underline{93.8}$ | $\underline{91.3}$ | | | |
| *Proprietary Reasoning LLMs* | | | | | | | | | | | | | | | | | | | | | |
| Claude 3.7 Sonnet [6] | $2.0_{40.0}$ | $2.1_{44.8}$ | $1.9_{35.6}$ | 42.0 | 46.9 | 37.5 | 49.0 | 36.5 | 60.6 | 4.0 | 3.1 | 4.8 | 93.5 | 95.8 | 91.3 | 93.0 | 90.6 | 95.2 | | | |
| Gemini 2.5 Flash [21] | $4.5_{1.0}$ | $3.1_{1.1}$ | $5.8_{0.9}$ | 5.5 | 4.2 | 6.7 | 88.0 | 84.4 | 91.3 | 13.5 | 7.3 | 19.2 | **100.0** | **100.0** | **100.0** | **100.0** | **100.0** | **100.0** | | | |
| Grok 3 mini [67] | $6.0_{65.5}$ | $4.2_{68.7}$ | $7.7_{62.5}$ | $\underline{\textbf{71.5}}$ | $\underline{72.9}$ | $\underline{\textbf{70.2}}$ | 24.0 | 16.7 | 30.8 | 19.5 | 11.5 | 26.9 | 53.5 | 63.5 | 44.2 | 91.0 | 94.8 | 87.5 | | | |
| Gemini 2.5 Pro [22] | $6.0_{1.0}$ | $7.3_{1.0}$ | $4.8_{1.0}$ | 7.0 | 8.3 | 5.8 | 88.5 | 83.3 | 93.3 | 19.0 | 12.5 | 25.0 | **100.0** | **100.0** | **100.0** | 99.5 | **100.0** | 99.0 | | | |
| o1 [45] | $8.0_{54.5}$ | $7.3_{55.2}$ | $8.7_{53.8}$ | 62.5 | 62.5 | 62.5 | 34.5 | 37.5 | 31.7 | 17.5 | 12.5 | 22.1 | 86.5 | 99.0 | 75.0 | 99.5 | **100.0** | 99.0 | | | |
| o3-mini [47] | $9.5_{53.0}$ | $7.3_{62.5}$ | $11.5_{44.3}$ | 62.5 | 69.8 | 55.8 | 37.0 | 34.4 | 39.4 | 22.0 | 17.7 | 26.0 | 77.5 | 92.7 | 63.5 | 95.0 | 96.9 | 93.3 | | | |
| o4-mini [48] | $15.5_{49.5}$ | $14.6_{48.9}$ | $16.3_{50.0}$ | 65.0 | 63.5 | 66.3 | 62.0 | 58.3 | 65.4 | 26.0 | 25.0 | 26.9 | 86.5 | 90.6 | 82.7 | 93.0 | 92.7 | 93.3 | | | |
| o3 [48] | $21.0_{16.0}$ | $18.8_{11.4}$ | $23.1_{20.2}$ | 37.0 | 30.2 | 43.3 | $\underline{93.5}$ | $\underline{91.7}$ | $\underline{95.2}$ | $\underline{39.5}$ | $\underline{28.1}$ | $\underline{50.0}$ | 91.5 | 99.0 | 84.6 | 97.0 | 96.9 | 97.1 | | | |
| **Average Accuracy (↑)** | $5.0_{38.0}$ | $4.5_{\downarrow42.9}$ | $5.5_{33.5}$ | 43.0 | 47.4 | 39.0 | 40.3 | 34.8 | 45.5 | 15.0 | 11.0 | 18.7 | 73.1 | 77.9 | 68.8 | 93.2 | 93.7 | 92.8 | | | |
| **Average Error Rate (↓)** | $95.0_{\uparrow38.0}$ | $95.5_{\uparrow42.9}$ | $94.5_{\uparrow33.5}$ | 57.0 | 52.6 | 61.0 | 59.7 | 65.2 | 54.5 | 85.0 | 89.0 | 81.3 | 26.9 | 22.1 | 31.2 | 6.8 | 6.3 | 7.2 | | | |

Table 4: Evaluation performance of chat and reasoning LLMs on the IneqMath benchmark (the test set). *Bnd.* denotes bound problems and *Rel.* denotes relation ones. We report: (1) *Overall Acc*, which reflects the correctness of both the final answer and intermediate steps; (2) *Answer Acc*, which measures final answer correctness alone; and (3) *Step Acc*, which evaluates the accuracy of intermediate steps across four error categories—*Toy Case*, *Logical Gap*, *Numerical Approximation*, and *Numerical Computation*. Blue superscripts ↓ indicate accuracy drop (*Overall Acc - Answer Acc*) from step-wise errors. Underlining denotes best result within each model category; **boldface** highlights best overall performance. Default max token limit for reasoning LLMs is 10K.

## 5.2 Main Evaluation Results

Table 4 presents the performance of the evaluated LLMs on IneqMath. Our analysis reveals several critical insights into current LLM capabilities for inequality proving:

**1) Reasoning LLMs achieve higher final-answer accuracy.** Models like o1 (62.5% *Answer Acc*) and Grok 3 mini (71.5% *Answer Acc*) significantly outperform their general-purpose chat counterparts (e.g., GPT-4o at 37.5%, Grok 3 at 54.5%) in identifying the correct final answer. This suggests that specialized architectures or training techniques improve their search ability to find final answers.

**2) Step-wise scrutiny reveals a dramatic performance drop.** The advantage in *Answer Acc* often masks underlying reasoning flaws. *Overall Acc* plummets when steps are evaluated. For instance, Grok 3 mini's accuracy drops by 65.5% (from 71.5% *Answer Acc* to 6.0% *Overall Acc*), and o3-mini by 53.0%. This stark discrepancy underscores the fragility of LLM-generated deductive chains.

**3) Robust proof construction remains a major challenge.** Even top models like o1 achieve low *Overall Acc* (8.0%). Many large models, despite moderate Answer Acc, also score poorly (e.g., Grok 3 at 3.5% *Overall Acc*). This indicates a fundamental gap between finding a plausible answer and constructing a mathematically rigorous, step-by-step derivation.

## 5.3 In-depth Study

**Failure solution analysis.** As shown in Table 4, where we report average error rates for overall accuracy, final-answer accuracy, and step-wise accuracy across four categories, the most common step-wise errors in LLM-generated solutions are logical gaps (85.0% average failure rate across models) and unjustified generalization from toy cases (59.7%). Less frequent, but still significant, are errors from numerical approximations (26.9%) and miscalculations (6.8%). A detailed inspection of incorrect solutions (see examples in §C.2.1-§C.2.4) highlights these prevalent error patterns, which often undermine proofs even when LLMs produce the correct final answer. Beyond these step-wise
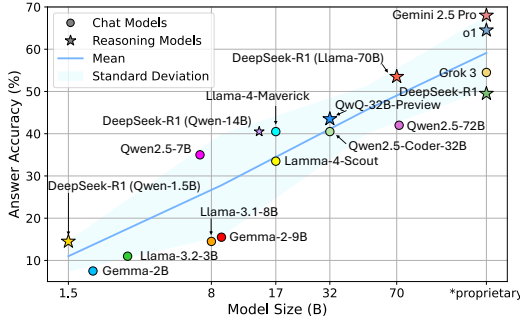
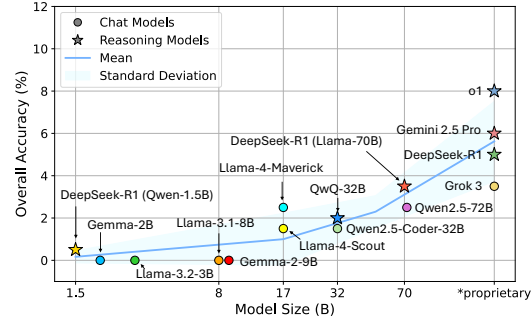Figure 6: Model-size scaling law (*Answer Acc*).



Figure 7: Model-size scaling law (*Overall Acc*).

errors, LLMs also struggle to derive correct final answers on complex problems (§C.2.5), indicating deeper challenges in theorem application and symbolic manipulation.

**Scaling law in model size.** Figure 6 shows how *final-answer* accuracy (which evaluates only the correctness of the final predicted answer) scales with model size for LLMs. As model size increases, we observe a steady improvement in answer accuracy, reflecting an empirical scaling law that larger models are better at inferring correct bounds and inequality relationships. However, this trend does not hold well when considering *overall accuracy*—which requires both a correct answer and valid intermediate reasoning steps—as shown in Figure 7. In this latter case, the scaling curve flattens, indicating that increased model size alone is insufficient to eliminate step-by-step reasoning errors.

**Scaling law in test-time computation.** Extended test-time computation, allowing longer reasoning chains, is a common strategy for complex problem-solving [14]. We investigated its impact on *overall accuracy* in INEQMATH by varying the maximum completion tokens for reasoning LLMs. Figure 8 shows that while models like Gemini 2.5 Pro and o3 initially improve with more tokens, performance gains saturate (e.g., beyond 20K tokens). This indicates that merely increasing computational budget offers diminishing returns for achieving rigorous, step-wise correct proofs, highlighting the need for more than just longer thought processes.



Figure 8: Scaling law in test-time computation for reasoning LLMs.

### 5.4 Exploring Improvement Strategies

**Retrieving relevant theorems as hints.** To assess theorem-based hints, we provide models with the top-$k$ most frequent theorems from our INEQMATH training corpus when solving a 40-problem test subset. As shown in Figure 9, providing one or two such theorems decreases *overall accuracy* for weaker models (e.g., Grok 3 mini, o3-mini, o4-mini), likely due to misapplication or distraction by potentially irrelevant information. Conversely, stronger models like Gemini 2.5 Pro benefit from these hints, suggesting advanced reasoning is crucial to effectively use such guidance. These results underscore the potential of theorem-guided reasoning but also highlight the critical need for more sophisticated theorem-retrieval mechanisms (e.g., RAG [28, 24]) to reliably enhance LLM performance in inequality proving. Detailed experiments are available in §C.4.



Figure 9: Model performance with retrieved theorems as hints.

**Self-improvement via critic as feedback.** Allowing an LLM to critique and revise its own reasoning has been shown to improve performance on complex tasks [78, 57]. To explore whether this holds for inequality proving, we randomly sampled 40 test problems from INEQMATH and ran one round of self-critique. As Figure 10 shows, self-critique consistently improves performance—e.g., Gemini 2.5 Pro's overall accuracy rises from 43% to 48%. This upward trend underscores self-critique as a promising, supervision-free method to enhance logical rigor and solution quality of LLMs in inequality reasoning. More details are in §C.5.
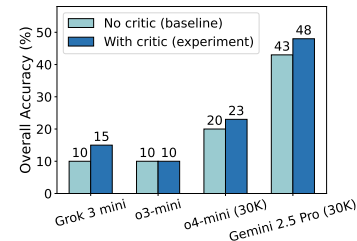


Figure 10: Model performance via self-critic as feedback.

# 6 Related Work

**Datasets for inequality and theorem proving.** One of the major bottlenecks in advancing LLM capabilities for inequality proving is the scarcity of suitable datasets. Existing resources fall short in several ways: general ATP collections like MiniF2F [82] and ProofNet [7] contain few inequalities; synthetic datasets such as INT [64] and AIPS [63] offer scale but often lack structural diversity due to their template-based generation; and curated collections like ChenNEQ [8] are often too small for extensive training. More fundamentally, most existing datasets [80, 59, 73, 29, 58, 26] adopt a *fully formal* representation, where problems and proofs are encoded in systems such as Lean [11] or Isabelle [42]. While formal mathematical reasoning offers correctness guarantees and is a vital research direction, LLMs, trained on vast corpora of natural language, often exhibit strong informal reasoning capabilities. Therefore, our INEQMATH adopts an *informal* perspective, reformulating inequality proof problems into two verifiable subtasks—bound estimation and relation prediction. These problems within INEQMATH were crafted and reviewed by IMO-level medalist experts. Other informal reasoning datasets [49, 39, 23, 37] typically lack annotated solutions, theorem references, or corresponding training data. To address these gaps, INEQMATH introduces 1,252 inequality problems for training, each annotated with theorems relevant to its solution, which comprises up to four steps.

**Methods for inequality and theorem proving.** Proving inequalities is complex, requiring intuition to identify tight bounds, strategic use of theorems, and precise symbolic manipulation. Traditional automated theorem provers (ATPs) primarily operate within formal systems like Lean [11] or Isabelle [42], requiring problems and proofs to be encoded in specialized languages. Inspired by the mathematical reasoning capabilities of LLMs [81], a significant body of recent work has focused on integrating LLMs with these formal ATPs. These approaches often model theorem proving as a Markov Decision Process (MDP), training LLMs to select appropriate tactics and premises to construct proofs within the formal system [1, 9, 62, 16, 18, 26, 34, 61, 72]. For instance, systems like Goedel-Prover [35] leverage large Lean corpora to train models for tactic prediction, enabling end-to-end formal proof generation. Other methods incorporate tree-search techniques to navigate the search space of premises within formal frameworks [65, 31, 70, 71].

LLMs are trained on vast natural language corpora, giving them strengths in informal reasoning—closer to how humans solve problems. This reveals an opportunity for methods that harness these informal abilities. Our work departs from formal paradigms by introducing an *informal yet verifiable* framework for inequality proving, designed to benchmark and enhance LLM performance in human-like problem solving, while exploring improvements such as theorem-guided reasoning and self-refinement.

**LLM-as-judge for math problem solving.** Reliable evaluation of mathematical problem-solving necessitates assessing not only the correctness of the final answer but also the logical soundness of each reasoning step, a significant challenge for automated systems. Traditional methods are often inadequate: expert annotation is labor-intensive and unscalable for large-scale evaluation [49, 39], while automated techniques such as string matching or value equivalence overlook crucial step-by-step proof correctness [25, 23, 37, 38]. While LLMs have shown promise as evaluators (*LLM-as-judge*), their capacity for detailed, step-wise mathematical judgment is still developing. For instance, existing step-level LLM judges [68, 17] may assess general step validity but often lack the granularity to identify nuanced reasoning flaws. Similarly, frameworks like EIC-Math [30] provide broad error categories but can miss subtle yet critical issues in multi-step derivations. To address these limitations and assess informal mathematical proofs like inequality solving, our *LLM-as-judge* framework combines a *final-answer* judge with four *step-wise* judges targeting common errors: toy case overgeneralization, logical gaps, unjustified numeric approximations, and numeric calculation mistakes.

# 7 Conclusion

In summary, we introduce an informal yet verifiable task formulation for inequality proving, decomposing it into bound estimation and relation prediction. Building on this, we release INEQMATH, an expert-curated benchmark of Olympiad-level inequalities with a training corpus featuring step-wise solutions and theorem annotations. Our novel *LLM-as-judge* evaluation framework, comprising a final-answer judge and four step-wise judges, enables a rigorous assessment. Our comprehensive evaluation of diverse leading LLMs reveals a critical gap: while LLMs may achieve high final-answer accuracy, this often plummets by up to 65.5% under step-wise scrutiny, with top models like o1 achieving less than 10% overall accuracy. This discrepancy exposes fragile deductive chains for current LLMs in constructing rigorous proofs. We further find that scaling model size or increasing test-time computation yields limited gains in overall proof correctness. Instead, our findings highlight promising research directions such as theorem-guided reasoning and self-refinement.

## Acknowledgments

## References

[1] Tudor Achim, Alex Best, Kevin Der, Mathïs Fédérico, Sergei Gukov, Daniel Halpern-Leister, Kirsten Henningsgard, Yury Kudryashov, Alexander Meiburg, Martin Michelsen, et al. Aristotle: Imo-level automated theorem proving. *arXiv preprint arXiv:2510.01346*, 2025.

[2] Meta AI. Llama 3.1 8b. https://huggingface.co/meta-llama/Llama-3.1-8B, 2024. Accessed: 2025-05-15.

[3] Meta AI. Llama 3.2 3b. https://huggingface.co/meta-llama/Llama-3.2-3B, 2024. Accessed: 2025-05-15.

[4] AI-MO / Project Numina & Kimi teams. Kimina-prover-distill-8b. https://huggingface.co/AI-MO/Kimina-Prover-Distill-8B, 2025.

[5] Alibaba Qwen Team. QwQ-32B. https://huggingface.co/Qwen/QwQ-32B, 2025. Hugging Face model card.

[6] Anthropic. Claude 3.7 Sonnet. https://www.anthropic.com/claude/sonnet, 2025. Anthropic model card.

[7] Zhangir Azerbayev, Bartosz Piotrowski, Hailey Schoelkopf, Edward William Ayers, and Dragomir Radev. Proofnet: Autoformalizing and formally proving undergraduate-level mathematics, 2024.

[8] Evan Chen. A brief introduction to olympiad inequalities, 2014. Accessed: 2025-03-19.

[9] Luoxin Chen, Jinming Gu, Liankai Huang, Wenhao Huang, Zhicheng Jiang, Allan Jie, Xiaoran Jin, Xing Jin, Chenggang Li, Kaijing Ma, et al. Seed-prover: Deep and broad reasoning for automated theorem proving. *arXiv preprint arXiv:2507.23726*, 2025.

[10] Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. TheoremQA: A theorem-driven question answering dataset. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7889–7901, Singapore, December 2023. Association for Computational Linguistics.

[11] Leonardo Mendonça de Moura, Soonho Kong, Jeremy Avigad, Floris van Doorn, and Jakob von Raumer. The lean theorem prover (system description). In *CADE*, 2015.

[12] DeepSeek-AI. DeepSeek-R1-Distill-Llama-70B. https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-70B, 2025. Hugging Face model card.

[13] DeepSeek-AI. DeepSeek-R1-Distill-Qwen-14B. https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-14B, 2025. Hugging Face model card.

[14] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.

[15] Kefan Dong and Tengyu Ma. Beyond limited data: Self-play llm theorem provers with iterative conjecturing and proving. *arXiv preprint arXiv:2502.00212*, 2025.

[16] Kefan Dong, Arvind V. Mahankali, and Tengyu Ma. Formal theorem proving by rewarding LLMs to decompose proofs hierarchically. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*, 2024.

[17] 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. Omni-MATH: A universal olympiad level mathematic benchmark for large language models. In *The Thirteenth International Conference on Learning Representations*, 2025.

[18] Fabian Gloeckle, Jannis Limperg, Gabriel Synnaeve, and Amaury Hayat. ABEL: Sample efficient online reinforcement learning for neural theorem proving. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*, 2024.

[19] Google DeepMind. Gemini 2.0 Flash. https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-0-flash, 2025. Vertex AI model card.

[20] Google DeepMind. Gemini 2.0 Flash-Lite. https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-0-flash-lite, 2025. Vertex AI model card.

[21] Google DeepMind. Gemini 2.5 Flash. https://cloud.google.com/vertex-ai/generative-ai/docs/models/gemini/2-5-flash, 2025. Vertex AI model card.

[22] Google DeepMind. Gemini 2.5 Pro. https://deepmind.google/technologies/gemini/pro/, 2025. Google DeepMind model card.

[23] Aryan Gulati, Brando Miranda, Eric Chen, Emily Xia, Kai Fronsdal, Bruno de Moraes Dumont, and Sanmi Koyejo. Putnam-AXIOM: A functional and static benchmark for measuring higher level mathematical reasoning. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*, 2024.

[24] Shailja Gupta, Rajesh Ranjan, and Surya Narayan Singh. A comprehensive survey of retrieval-augmented generation (rag): Evolution, current landscape and future directions. *arXiv preprint arXiv:2410.12837*, 2024.

[25] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

[26] Jiewen Hu, Thomas Zhu, and Sean Welleck. miniCTX: Neural theorem proving with (long-)contexts. In *The Thirteenth International Conference on Learning Representations*, 2025.

[27] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.

[28] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc., 2020.

[29] Wenda Li, Lei Yu, Yuhuai Wu, and Lawrence C Paulson. Isarstep: a benchmark for high-level mathematical reasoning. In *International Conference on Learning Representations*, 2021.

[30] Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. Evaluating mathematical reasoning of large language models: A focus on error identification and correction. *arXiv preprint arXiv:2406.00755*, 2024.

[31] Yang Li, Dong Du, Linfeng Song, Chen Li, Weikang Wang, Tao Yang, and Haitao Mi. Hunyuanprover: A scalable data synthesis framework and guided tree search for automated theorem proving. *arXiv preprint arXiv:2412.20735*, 2024.

[32] Zenan Li, Zhaoyu Li, Wen Tang, Xian Zhang, Yuan Yao, Xujie Si, Fan Yang, Kaiyu Yang, and Xiaoxing Ma. Proving olympiad inequalities by synergizing LLMs and symbolic reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025.

[33] Zhenwen Liang, Linfeng Song, Yang Li, Tao Yang, Feng Zhang, Haitao Mi, and Dong Yu. Mps-prover: Advancing stepwise theorem proving by multi-perspective search and data curation. *arXiv preprint arXiv:2505.10962*, 2025.

[34] Haohan Lin, Zhiqing Sun, Sean Welleck, and Yiming Yang. Lean-STar: Learning to interleave thinking and proving. In *The Thirteenth International Conference on Learning Representations*, 2025.

[35] Yong Lin, Shange Tang, Bohan Lyu, Jiayun Wu, Hongzhou Lin, Kaiyu Yang, Jia Li, Mengzhou Xia, Danqi Chen, Sanjeev Arora, et al. Goedel-prover: A frontier model for open-source automated theorem proving. *arXiv preprint arXiv:2502.07640*, 2025.

[36] Yong Lin, Shange Tang, Bohan Lyu, Ziran Yang, Jui-Hui Chung, Haoyu Zhao, Lai Jiang, Yihan Geng, Jiawei Ge, Jingruo Sun, Jiayun Wu, Jiri Gesi, Ximing Lu, David Acuna, Kaiyu Yang, Hongzhou Lin, Yejin Choi, Danqi Chen, Sanjeev Arora, and Chi Jin. Goedel-prover-v2: Scaling formal theorem proving with scaffolded data synthesis and self-correction, 2025.

[37] Junnan Liu, Hongwei Liu, Linchen Xiao, Ziyi Wang, Kuikun Liu, Songyang Gao, Wenwei Zhang, Songyang Zhang, and Kai Chen. Are your llms capable of stable reasoning? *arXiv preprint arXiv:2412.13147*, 2024.

[38] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*, 2024.

[39] Yujun Mao, Yoon Kim, and Yilun Zhou. CHAMP: A competition-level dataset for fine-grained analyses of LLMs' mathematical reasoning capabilities. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13256–13274, Bangkok, Thailand, August 2024. Association for Computational Linguistics.

[40] Meta Platforms, Inc. Llama-4-Maverick-17B-128E. https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E, 2025. Hugging Face model card.

[41] Meta Platforms, Inc. Llama-4-Scout-17B-16E. https://huggingface.co/meta-llama/Llama-4-Scout-17B-16E, 2025. Hugging Face model card; accessed 2025-05-12.

[42] Tobias Nipkow, Markus Wenzel, and Lawrence C. Paulson. *Isabelle/HOL: A Proof Assistant for Higher-Order Logic*. Springer, 2002.

[43] OpenAI. GPT-4o. https://openai.com/index/hello-gpt-4o/, 2024. OpenAI blog post.

[44] OpenAI. GPT-4o mini. https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/, 2024. OpenAI blog post.

[45] OpenAI. OpenAI o1. https://openai.com/o1/, 2024. OpenAI official announcement.

[46] OpenAI. GPT-4.1. https://openai.com/index/gpt-4-1/, 2025. OpenAI model announcement.

[47] OpenAI. Openai o3-mini system card, January 2025. Accessed: 2025-03-19.

[48] OpenAI. OpenAI o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/, 2025. OpenAI official announcement.

[49] Ivo Petrov, Jasper Dekoninck, Lyuben Baltadzhiev, Maria Drencheva, Kristian Minchev, Mislav Balunović, Nikola Jovanović, and Martin Vechev. Proof or bluff? evaluating llms on 2025 usa math olympiad. *arXiv preprint arXiv:2503.21934*, 2025.

[50] Gabriel Poesia, David Broman, Nick Haber, and Noah Goodman. Learning formal mathematics from intrinsic motivation. *Advances in Neural Information Processing Systems*, 37:43032–43057, 2024.

[51] Qwen Team. Qwen2.5-72B. https://huggingface.co/Qwen/Qwen2.5-72B, 2024. Hugging Face model card.

[52] Qwen Team. Qwen2.5-7B. https://huggingface.co/Qwen/Qwen2.5-7B, 2024. Hugging Face model card.

[53] Qwen Team. Qwen3-235B-A22B. https://huggingface.co/Qwen/Qwen3-235B-A22B, 2025. Hugging Face model card.

[54] Z. Z. Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanjia Zhao, Liyue Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, Z. F. Wu, Zhibin Gou, Shirong Ma, Hongxuan Tang, Yuxuan Liu, Wenjun Gao, Daya Guo, and Chong Ruan. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition, April 2025.

[55] Z. Z. Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanjia Zhao, Liyue Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, Z. F. Wu, Zhibin Gou, Shirong Ma, Hongxuan Tang, Yuxuan Liu, Wenjun Gao, Daya Guo, and Chong Ruan. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition, 2025.

[56] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.

[57] Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Lei Han, Haitao Mi, and Dong Yu. Toward self-improvement of LLMs via imagination, searching, and criticizing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

[58] George Tsoukalas, Jasper Lee, John Jennings, Jimmy Xin, Michelle Ding, Michael Jennings, Amitayush Thakur, and Swarat Chaudhuri. Putnambench: Evaluating neural theorem-provers on the putnam mathematical competition. *arXiv preprint arXiv:2407.11214*, 2024.

[59] Nguyen Duy Tung. 567 nice and hard inequalities, April 2012. Accessed: 2025-03-19.

[60] Haiming Wang, Mert Unsal, Xiaohan Lin, Mantas Baksys, Junqi Liu, Marco Dos Santos, Flood Sung, Marina Vinyes, Zhenzhe Ying, Zekai Zhu, Jianqiao Lu, Hugues de Saxc'e, Bolton Bailey, Chendong Song, Chenjun Xiao, Dehao Zhang, Ebony Zhang, Frederick Pu, Han Zhu, Jiawei Liu, Jonas Bayer, Julien Michel, Longhui Yu, Léo Dreyfus-Schmidt, Lewis Tunstall, Luigi Pagani, Moreira Machado, Pauline Bourigault, Ran Wang, Stanislas Polu, Thibaut Barroyer, Wen-Ding Li, Yazhe Niu, Yann Fleureau, Yangyang Hu, Zhouliang Yu, Zihan Wang, Zhilin Yang, Zhengying Liu, and Jia Li. Kimina-prover preview: Towards large formal reasoning models with reinforcement learning. *arXiv preprint arXiv:2504.11354*, 2025.

[61] Haiming Wang, Huajian Xin, Chuanyang Zheng, Zhengying Liu, Qingxing Cao, Yinya Huang, Jing Xiong, Han Shi, Enze Xie, Jian Yin, Zhenguo Li, and Xiaodan Liang. LEGO-prover: Neural theorem proving with growing libraries. In *The Twelfth International Conference on Learning Representations*, 2024.

[62] Ruida WANG, Rui Pan, Yuxin Li, Jipeng Zhang, Yizhen Jia, Shizhe Diao, Renjie Pi, Junjie Hu, and Tong Zhang. MA-lot: Model-collaboration lean-based long chain-of-thought reasoning enhances formal theorem proving. In *Forty-second International Conference on Machine Learning*, 2025.

[63] Chenrui Wei, Mengzhou Sun, and Wei Wang. Proving olympiad algebraic inequalities without human demonstrations. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.

[64] Yuhuai Wu, Albert Jiang, Jimmy Ba, and Roger Baker Grosse. {INT}: An inequality benchmark for evaluating generalization in theorem proving. In *International Conference on Learning Representations*, 2021.

[65] Zijian Wu, Suozhi Huang, Zhejian Zhou, Huaiyuan Ying, Jiayu Wang, Dahua Lin, and Kai Chen. Internlm2. 5-stepprover: Advancing automated theorem proving via expert iteration on large-scale lean problems. *arXiv preprint arXiv:2410.15700*, 2024.

[66] xAI. Grok 3. https://x.ai/news/grok-3, 2025. xAI official announcement.

[67] xAI. Grok 3 Mini. https://x.ai/news/grok-3, 2025. xAI official announcement.

[68] Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. Evaluating mathematical reasoning beyond accuracy. *CoRR*, abs/2404.05692, 2024.

[69] Huajian Xin, Daya Guo, Zhihong Shao, Z.Z. Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. Advancing theorem proving in LLMs through large-scale synthetic data. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*, 2024.

[70] Ran Xin, Chenguang Xi, Jie Yang, Feng Chen, Hang Wu, Xia Xiao, Yifan Sun, Shen Zheng, and Kai Shen. Bfs-prover: Scalable best-first tree search for llm-based automatic theorem proving. *arXiv preprint arXiv:2502.03438*, 2025.

[71] Ran Xin, Zeyu Zheng, Yanchen Nie, Kun Yuan, and Xia Xiao. Scaling up multi-turn off-policy rl and multi-agent tree search for llm step-provers. *arXiv preprint arXiv:2509.06493*, 2025.

[72] Yu Xuejun, Jianyuan Zhong, Zijin Feng, Pengyi Zhai, Roozbeh Yousefzadeh, Wei Chong Ng, Haoxiong Liu, Ziyi Shou, Jing Xiong, Yudong Zhou, et al. Mathesis: Towards formal theorem proving from natural languages. *arXiv preprint arXiv:2506.07047*, 2025.

[73] Kaiyu Yang and Jia Deng. Learning to prove theorems via interacting with proof assistants. In *International Conference on Machine Learning*, pages 6984–6994. PMLR, 2019.

[74] Kaiyu Yang, Aidan M Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. Leandojo: Theorem proving with retrieval-augmented language models. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

[75] Ziyu Ye, Jiacheng Chen, Jonathan Light, Yifei Wang, Jiankai Sun, Guohao Li, Mac Schwager, Philip Torr, Yuxin Chen, Kaiyu Yang, Yisong Yue, and Ziniu Hu. Reasoning in reasoning: A hierarchical framework for (better and faster) neural theorem proving. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*, 2024.

[76] Huaiyuan Ying, Zijian Wu, Yihan Geng, JIayu Wang, Dahua Lin, and Kai Chen. Lean workbook: A large-scale lean problem set formalized from natural language math problems. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.

[77] Zhouliang Yu, Ruotian Peng, Keyi Ding, Yizhe Li, Zhongyuan Peng, Minghao Liu, Yifan Zhang, Zheng Yuan, Huajian Xin, Wenhao Huang, Yandong Wen, Ge Zhang, and Weiyang Liu. Formalmath: Benchmarking formal mathematical reasoning of large language models, 2025.

[78] Mert Yuksekgonul, Federico Bianchi, Jack Boen, et al. Optimizing generative ai by back-propagating language model feedback. *Nature*, 639:609–616, 2025. Published 19 March 2025.

[79] Ziyin Zhang, Jiahao Xu, Zhiwei He, Tian Liang, Qiuzhi Liu, Yansi Li, Linfeng Song, Zhengwen Liang, Zhuosheng Zhang, Rui Wang, et al. Deeptheorem: Advancing llm reasoning for theorem proving through natural language and reinforcement learning. *arXiv preprint arXiv:2505.23754*, 2025.

[80] Haoyu Zhao, Yihan Geng, Shange Tang, Yong Lin, Bohan Lyu, Hongzhou Lin, Chi Jin, and Sanjeev Arora. Ineq-Comp: Benchmarking human-intuitive compositional reasoning in automated theorem proving on inequalities. *arXiv preprint arXiv:2505.12680*, 2025.

[81] Zhiyu Zhao, Yongcheng Zeng, Ning Yang, and Guoqing Liu. Enhancing mathematical reasoning in language models through focused differentiation training, 2025.

[82] Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. minif2f: a cross-system benchmark for formal olympiad-level mathematics. In *International Conference on Learning Representations*, 2022.

# Supplementary Materials for
# Solving Inequality Proofs with Large Language Models

## Appendix Contents

# A Dataset Curation Details

## A.1 Training data curation.

Training problems were sourced from two advanced textbooks featuring graduate-level and Olympiad-style inequality proof problems. We parsed these textbooks to extract proof problems, their step-wise solutions, and relevant theorems. We developed two LLM-based rephrasers to transform each source problem into two sub-tasks defined in §2: bound estimation and relation prediction. For instance, a source problem like "Prove $a + b \geq 2\sqrt{ab}$ for $\forall a, b \in \mathbb{R}^+$" would be rephrased into a bound estimation task (e.g., "Determine the maximal constant $C$ such that $a + b \geq C\sqrt{ab}$ holds for $\forall a, b \in \mathbb{R}^+$") and a relation prediction task (e.g., "Determine the inequality relation in the expression $a + b$ ( ) $2\sqrt{ab}$ that holds for $\forall a, b \in \mathbb{R}^+$").

Crucially, while rephrased problems are altered from the source proof problem in the format, they preserve the core mathematical reasoning and solution steps—such as applying relevant theorems, determining boundary conditions, and verifying inequalities. An annotation tool (see §A.2) was developed to facilitate human expert review and correction of the LLM-rephrased problems. Extracted theorems were curated, each including its name, a natural language definition, and a list of training problems where it is applicable.

## A.2 Data Annotation Tool



Figure 11: The interface of our developed tool for checking and editing the bound problems.

Figure 12: The interface of our developed tool for checking and editing the relation problems.

## A.3 Prompts for Rephrasing Problems

---

### Prompt for Rephrasing Proofs to Bound Problems

**Task:** Transform the given inequality problem into a bound prediction problem by introducing a constant C and determining its optimal value.

**Instructions:**
1. Analyze the original problem, focusing on its structure and potential for transformation.
2. Introduce a constant $C$ by either replacing an existing constant or creating a new relationship between expressions.
3. Determine whether to find the minimal or maximal value of $C$ that satisfies the inequality for all relevant variables.
4. Consider factors such as homogeneity, existing constraints, and the domain of variables (e.g., positive reals, all reals).
5. Ensure the rephrased problem maintains the mathematical essence and constraints of the original.

**Output Format:**
Provide your response in the following structure:
`<Analysis>`: Concise explanation of key features and transformation approach.
`<Conclusion>`: YES or NO, followed by a brief summary of the transformation.
`<Rephrased Problem>`: Transformed problem statement, focusing on finding the optimal $C$.
`<Answer>`: $C = < value >$.

**Key Considerations:**
1. For double inequalities, attempt to rephrase as a single bound prediction problem when possible.
2. In homogeneous inequalities, focus on the ratios between variables rather than their absolute values.
3. Incorporate any existing constraints into the rephrased version of the problem.
4. Clearly specify the domain of the variables in the rephrased problem statement.
5. Ensure that the rephrased problem is logically equivalent to the original.

**Example:**

---

Original problem: Let $a, b, c \in \mathbb{R}^+$. Prove the inequality

$$\frac{abc}{(1+a)(a+b)(b+c)(c+16)} \leq \frac{1}{81}$$

<Analysis>: To turn this into a bound prediction problem, we can focus on the following steps:
1. The left side is a rational expression that is always positive for $a, b, c \in \mathbb{R}^+$.
2. The right side is a fixed constant $\frac{1}{81}$.
3. We replace the constant $\frac{1}{81}$ with a variable $C$ and ask: What is the smallest $C$ such that the inequality holds for all positive $a, b, c$?
4. This approach allows us to determine the tightest possible upper bound for the left-hand expression.
5. If we find the smallest $C$ that works, we prove the original inequality and show it's the best possible.

<Conclusion>: YES, the inequality can be rephrased as a bound prediction problem. By replacing the constant $\frac{1}{81}$ with a variable $C$, we can determine the tightest upper bound for the given rational expression, effectively proving the original inequality and demonstrating its optimality.

<Rephrased problem>:
Coll Determine the minimal constant $C$ such that the following inequality holds for all $a, b, c$:

$$\frac{abc}{(1+a)(a+b)(b+c)(c+16)} \leq C.$$

<Answer>: $C = \frac{1}{81}$.

**Now, please rewrite the following problem:**
Original problem: {problem}

---

## Prompt for Rephrasing Proofs to Relation Problems

**Task:** Transform the given inequality proof problem into a relation prediction problem.

**Instructions:**
1. Analyze the original problem, identifying key components such as variables, domains, conditions, and the main inequality.
2. Rephrase the problem by maintaining the original expressions and replacing the relation symbol with a blank to be filled.
3. Preserve any additional conditions or constraints from the original problem in your rephrased version.
4. Change the task from "Prove" to "Determine the correct inequality relation to fill in the blank."
5. Provide a set of options for the relation, always including $\leq, \geq, =, <, >$, and "None of the above".
6. Determine the correct answer based on your modification and analysis.

**Output Format:**
Provide your response in the following structure:

<Analysis>: Detailed step-by-step analysis of the original problem and your approach to rephrasing it.

<Conclusion>: YES or NO, followed by a brief explanation of whether and how the problem can be effectively rephrased.

<Rephrased Problem>:
  Transformed problem statement.
  Options: (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

<Answer>: Option letter.

**Key Considerations:**
1. Maintain the original mathematical expressions and any given conditions as much as possible.
2. Ensure the rephrased problem captures the essence and complexity of the original problem.
3. For problems with multiple inequalities, focus on one main inequality for the relation prediction task.
4. When dealing with complex fractions or expressions, keep them intact to maintain the problem's difficulty level.
5. If the relation depends on specific values of the variables or cannot be definitively determined, consider using "None of the above" as the correct answer.

**Example:**
Original problem: Let $a, b, c \in \mathbb{R}^+$. Prove the inequality

$$\frac{abc}{(1+a)(a+b)(b+c)(c+16)} \leq \frac{1}{81}$$

`<Analysis>`: To rephrase it to a relation prediction problem, we can focus on the following steps:

1. The original problem is a proof task for an inequality involving positive real numbers $a$, $b$, and $c$.
2. The left side of the inequality is a complex fraction $\frac{abc}{(1+a)(a+b)(b+c)(c+16)}$.
3. The right side is a constant fraction $\frac{1}{81}$.
4. The original inequality uses the "less than or equal to" ($\leq$) relation, which needs to hold for all positive real values of $a$, $b$, and $c$.
5. We can transform the proof task into determining the correct relation between the left and right sides of the inequality.

`<Conclusion>`: YES, the inequality can be effectively rephrased as a relation prediction problem.

`<Rephrased problem>`:
Let $a, b, c \in \mathbb{R}^+$. Consider the following inequality:

$$\frac{abc}{(1+a)(a+b)(b+c)(c+16)} \quad (\ ) \quad \frac{1}{81}.$$

Determine the correct inequality relation to fill in the blank.
Options: (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

`<Answer>`: A

**Now, please rewrite the following problem:**
Original problem: {problem}

## A.4 Benchmark Examples

---

**INEQMATH Testing Example 1: Bound Problem**

**Question:** Let $x, y, z > 0$ such that $x + y + z = 1$. Determine the minimal constant $C$ such that the following inequality holds for all $x, y, z$:

$$xy(y + 4z) + yz(z + 4x) + zx(x + 4y) \leq C.$$

---

**INEQMATH Testing Example 2: Bound Problem**

**Question:** Let $a_1, a_2, \ldots, a_n$ be real numbers and $S$ be a non-empty subset of $\{1, 2, \ldots, n\}$. Find the largest constant $C$ such that the following inequality holds for all $a_1, a_2, \ldots, a_n$ and $S$:

$$2C \left( \sum_{i \in S} a_i \right)^2 \leq \sum_{1 \leq i \leq j \leq n} (a_i + \cdots + a_j)^2.$$

---

**INEQMATH Testing Example 3: Bound Problem**

**Question:** Let $a_1, a_2, \ldots, a_n > 0$ such that $a_1 + a_2 + \ldots + a_n < 1$. Determine the minimal constant $C$ such that the following inequality holds for all $a_1, a_2, \ldots, a_n$:

$$\frac{a_1 \cdot a_2 \ldots a_n \left(1 - a_1 - a_2 - \ldots - a_n\right)}{\left(a_1 + a_2 + \ldots + a_n\right)\left(1 - a_1\right)\left(1 - a_2\right)\ldots\left(1 - a_n\right)} \leq C \frac{3}{n^{n-1}}.$$

**INEQMATH Testing Example 4: Relation Problem**

**Question:** Let $a, b, c, x, y, z \in \mathbb{R}$ be real numbers such that $a + b + c = 1$ and $x^2 + y^2 + z^2 = 1$. Consider the following expression:

$$a(x + b) + b(y + c) + c(z + a) \quad (\quad) \quad 1.$$

Determine the correct inequality relation to fill in the blank.

**Options:** (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

---

**INEQMATH Testing Example 5: Relation Problem**

**Question:** In the plane of the acute-angled triangle $\triangle ABC$, let $L$ be a line such that $u, v, w$ are the lengths of the perpendiculars from $A, B, C$ respectively to $L$. Consider the following inequality:

$$u^2 \tan A + v^2 \tan B + w^2 \tan C \quad (\quad) \quad \Delta.$$

where $\Delta$ is the area of the triangle. Determine the correct inequality relation to fill in the blank.

**Options:** (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

---

**INEQMATH Testing Example 6: Relation Problem**

**Question:** Let $a, b, c$ be the sides of any triangle. Consider the following inequality:

$$3 \left( \sum_{cyc} ab \left(1 + 2\cos(c)\right) \right) \quad (\quad) \quad 2 \left( \sum_{cyc} \sqrt{\left(c^2 + ab(1 + 2\cos(c))\right)\left(b^2 + ac(1 + \cos(b))\right)} \right).$$

Determine the correct inequality relation to fill in the blank.

**Options:** (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

# B Fine-grained Informal Judge Details

## B.1 Final Answer Judge

---

**Prompt for Final Answer Judge: Answer Extraction for Bound problems**

You are an expert in extracting numbers from answer sentences. Below are examples of sentences and the corresponding numbers:

Example 1: answer is $C = 2$.
Answer: $C = 2$

Example 2: answer is $C = \frac{1}{\sqrt{2}}$.
Answer: $C = \frac{1}{\sqrt{2}}$

Example 3: answer is $\boxed{C = 2}$.
Answer: $C = 2$

Now, extract the number from the following sentence: {`answer_sentence`}.

Make sure to return the answer in the format as "C=<`extracted_answer`>", where <`extracted_answer`> is the extracted number or expression.

---

**Prompt for Final Answer Judge: Answer Extraction for Relation Problems**

You are an expert in extracting option letters (A, B, C, D, E, F) from answer sentences.

The options are given below:
A: (A) $\leq$
B: (B) $\geq$
C: (C) $=$
D: (D) $<$
E: (E) $>$
F: (F) None of the above

Below are examples of sentences and the corresponding option letters:

Example 1: answer is (B) $\geq$.
Answer: B

Example 2: answer is (E) >.
Answer: E

Example 3: answer is: $\boxed{\leq}$.
Answer: A

Now, extract the option letter from the following sentence: {`answer_sentence`}.

Make sure to return the option letter only, without any other characters.

---

**Prompt for Final Answer Judge: Answer Equivalence Verification**

You are an expert in verifying mathematical expression equivalence. Analyze if two expressions are exactly equivalent by following these strict rules:

**Required Analysis Steps:**
1. Check if both expressions are valid mathematical forms.
2. If either expression is not mathematical (e.g., text or undefined), return `False`.
3. For numerical expressions:
   - Direct equality (e.g., $2 = 2$) → `True`.

---

- Different representations of same value (e.g., $1/2 = 0.5$, $\sqrt{4} = 2$) → `True`.
- Decimal approximations vs exact values (e.g., $2\pi \neq 6.28318$) → `False`.
4. For algebraic expressions:
   - Must have clear, valid transformation path between forms.
   - If transformation requires multiple non-obvious steps → `False`.
   - Verify equivalence through algebraic proof when possible.
   - For complex expressions, use techniques like squaring or substitution to verify.

**Equivalence Criteria:**
   - Must have exactly same deterministic value.
   - Must be provably equivalent through valid mathematical operations.
   - Different notations of same exact value are equivalent.
   - Decimal approximations are NOT equivalent to exact expressions.
   - No rounding or approximations allowed.
   - If equivalence cannot be conclusively proven → `False`.

**Example Pairs and their Analysis:**
Ground truth: $C = 2$
Prediction: $C = 2$
Analysis: The expressions are identical in both form and value, representing the same integer 2.
Equivalent: `True`

Ground truth: $C = 1.5$
Prediction: $C = \frac{3}{2}$
Analysis: The decimal 1.5 and fraction $\frac{3}{2}$ are different representations of the same number ($1.5 = \frac{3}{2}$).
Equivalent: `True`

Ground truth: $C = 2\pi$
Prediction: $C = 6.28318530718$
Analysis: While 6.28318530718 is a decimal approximation of $2\pi$, they are not symbolically equivalent expressions.
Equivalent: `False`

Ground truth: $C = \sqrt{\frac{1}{6}}$
Prediction: $C = \frac{1}{\sqrt{6}}$
Analysis: These are equivalent through the property $\sqrt{\frac{a}{b}} = \frac{\sqrt{a}}{\sqrt{b}}$ when $a, b > 0$.
Equivalent: `True`

Ground truth: $C = \sqrt{\frac{3}{2}}$
Prediction: $C = \frac{3}{2\sqrt{2}}$
Analysis: The expressions differ as proven when squared: $(\sqrt{\frac{3}{2}})^2 = \frac{3}{2} \neq \frac{9}{8} = (\frac{3}{2\sqrt{2}})^2$.
Equivalent: `False`

**Now analyze these expressions:**
Ground truth: {ground_truth}
Prediction: {prediction}

## B.2 Toy Case Judge

**Prompt for Toy Case Judge**

**Task:** Evaluate the logical rigor of a solution to an inequality problem, focusing specifically on whether the direction of the inequality was justified using toy cases or special value substitution.

**Instructions:**
1. Carefully read the reasoning process used to solve the inequality.
2. Identify whether the direction of the inequality was determined by testing special values, trying toy cases, or relying on extreme-case analysis, rather than providing a general proof valid over the entire domain.

3. If the model uses a toy case (e.g., setting a variable to 0, 1, or choosing symmetric/equal values) or considers a variable tending to $0$ or $\infty$ (extreme-case reasoning) to conclude the inequality direction, this should be flagged as logically unsound unless it is later supported by a rigorous or general argument.
4. Substituting special values for the purpose of verifying equality or testing sharpness is acceptable and should not be flagged.
5. If a toy case is used to show that the inequality does not hold (i.e., the two sides are incomparable), this is acceptable and should not be flagged.
6. Trying toy cases or substituting special values for the purpose of exploring or analyzing the problem—without using them to directly conclude the inequality direction—is acceptable and should not be flagged.
7. The goal is to confirm whether the final conclusion is justified for *all variables in the given domain* by using sound and formal reasoning.

**Output Format:**
<Analysis>: Brief explanation of whether toy cases, special values, or extreme-case reasoning were used to justify inequality direction, and whether this undermines the rigor of the proof.
<Flagged Reasoning Step (if applicable)>: Quote or summarize the step(s) where toy cases or extreme cases were used to justify the inequality.
<Answer>: True or False. It indicates whether the solution's reasoning is logically valid (True) or contains flaws due to unjustified use of toy cases, special value substitution, or extreme-case reasoning (False).

**Key Considerations:**
1. Check which steps are used to justify the final general conclusion.
2. Toy case examples should not be used to conclude general inequality directions.
3. These techniques are allowed for illustration or equality testing, but not as a replacement for a full justification.
4. Extreme-case analysis used as a substitute for proof to justify the final general conclusion should be flagged.
5. Do not provide improvement suggestions—simply judge whether the logic used is valid or flawed.
6. Only reasoning steps that fail to generalize and are used to justify the final general conclusion should be flagged.

**Examples of Inputs and Outputs:**
{examples}

**Now analyze the following problem and solution:**
Original Problem: {query}
Solution: {response}

---

**Judge Example 1: Toy Case**

**Response:**
Let's test the specific values.
`Case 1:` If $a = b = c$, then:
$$\left(\frac{a}{b+c}\right)^2 + \frac{b}{c+a} + \left(\frac{c}{a+b}\right)^2 = \frac{1}{4} + \frac{1}{2} + \frac{1}{4} = 1.$$
`Case 2:` If $a = 0$ and $b = c = 1$, then:
$$\left(\frac{a}{b+c}\right)^2 + \frac{b}{c+a} + \left(\frac{c}{a+b}\right)^2 = 0 + 1 + 1 = 2.$$
From the test cases, we can conclude that the minimum value of the expression was 1. Therefore, the expression is always at least 1. The answer is $\boxed{C = 1}$.

**Judge Output:**
<Analysis>: The minimum of the expression is derived from two special cases(i.e. $a = b = c$ and $a = 0$, $b = c = 1$). This is incorrect as it substitutes generality with selected toy cases.

---

24

> `<Flagged Reasoning Step (if applicable)>`: From the test cases, we can conclude that the minimum value of the expression is 1. Therefore, the expression is always at least 1.
>
> `<Answer>`: $\boxed{False}$

## B.3 Logical Gap Judge

**Prompt for Logical Gap Judge**

**Task:** Evaluate the logical rigor of a proposed solution to a mathematical inequality problem. Focus on whether the reasoning includes non-trivial claims made without justification, logical leaps, or unsupported assertions, while allowing for valid optimization-based, algebraic, or analytic analysis when properly demonstrated. This prompt does not evaluate whether the direction of an inequality was justified using toy cases, special values, or asymptotic behavior; that aspect is handled separately.

**Instructions:**
1. Carefully read the entire reasoning process used to solve the inequality.
2. Identify whether the solution includes:
   - Any non-obvious (non-trivial) claims or transformations without justification or explanation.
   - Any logical gaps or skipped steps that lead to intermediate or final conclusions.
3. All significant transformations—especially involving inequalities, bounds, or extremal behavior—must be supported by: algebraic manipulation, well-known identities or theorems, valid analytical tools (e.g., convexity, derivatives, limits) or step-by-step numeric or symbolic reasoning.
4. Optimization methods (e.g., Lagrange multipliers, derivative-based analysis) are valid only if the analysis is explicitly shown:
   - If a solution invokes optimization or analytical techniques, it must demonstrate key steps, derivative conditions, or critical point verification.
   - Statements such as "solving the constrained optimization problem confirms..." without any derivation or argument are considered unjustified.
   - You do not need to assess whether toy cases, special values, or extreme behavior were used to infer the inequality direction. That responsibility lies outside the scope of this Judge.
5. Simple algebra or widely known transformations (e.g., AM-GM, factoring identities) may be used without full derivation.
6. The goal is to assess whether each important conclusion within the reasoning—not just the final answer—is logically supported and rigorously justified.

**Output Format:**
`<Analysis>`: Step-by-step explanation of whether the reasoning is logically sound. Highlight any unjustified claims or skipped steps, unless they are supported by valid asymptotic, numeric, or analytic reasoning.
`<Flagged Reasoning Step (if applicable)>`: Quote or summarize the specific step(s) where an unjustified claim or logical leap occurred.
`<Answer>`: True or False. True if the reasoning is valid; False if it contains unjustified steps or unsupported claims.

**Key Considerations:**
1. Focus on whether each major step (not just the final answer) is logically justified.
2. Non-trivial algebraic identities or inequality steps must be explained unless they are well-known.
3. Minor simplifications and standard techniques are acceptable without proof.
4. Do not flag steps involving toy cases, extreme values, or special substitutions used to infer inequality direction—those are out of scope here.
5. Claims like "a numerical check shows" must include specific values, results, or graphs to be valid.
6. Optimization-based arguments (e.g., Lagrange multipliers, critical point methods) must include demonstrated steps or analytic structure. If only the method is named but not applied, the reasoning should be flagged.
7. Do not flag steps used solely for equality verification, sharpness testing, or illustration.
8. Do not provide improvement suggestions—simply judge whether the logic is valid or flawed as presented.

**Examples of Inputs and Outputs:**
`{examples}`

**Now analyze the following problem and solution:**
Original Problem: {query}
Solution: {response}

---

**Judge Example 2: Logical Gap**

**Response:**
...
A numerical check confirms that the minimum occurs when $x = 1$, and the expression equals 4. Therefore, the largest constant $C$ such that the inequality holds for all positive $a, b, c$ is:

$$\boxed{C = 4}$$

**Judge Output:**
<Analysis>: The statement "A numerical check confirms that the minimum occurs when $x = 1$" is made without showing any actual numerical data, evaluations, or visual/analytic evidence to support this claim.

<Flagged Reasoning Step (if applicable)>: A numerical check confirms that the minimum occurs when $x = 1$.

<Answer>: $\boxed{False}$

---

## B.4 Numerical Approximation Judge

**Prompt for Numerical Approximation Judge**

**Task:** Evaluate the logical rigor of a solution to an inequality problem, specifically checking whether approximate numerical substitutions (e.g., replacing $\sqrt{2}$ with $1.414$) were improperly used in the reasoning process.

**Instructions:**
1. Carefully read the entire reasoning process used to solve the inequality.
2. Identify whether the solution includes:
   - Any replacement of exact expressions (such as radicals, fractions, or constants like $\pi$) with approximate decimal values.
3. Strict rules for use of approximate values:
   - If approximated values are directly involved in any operations (such as addition, subtraction, multiplication, or division), it must immediately be considered invalid, regardless of whether the operation is for comparing sizes or for further reasoning!
   - Examples of invalid actions: Approximating $\sqrt{5} \approx 2.236$ and then using it to compute $\sqrt{5} + 3$ approximately, or Approximating $\pi \approx 3.14$ and then evaluating $\pi/2$ based on $3.14$.
4. Approximate substitutions are allowed only under the following conditions: If approximate numerical comparison is used between simple numbers (e.g., $\sqrt{4}$, $\frac{1}{2}$, $\sqrt{2}$) that humans can readily estimate, it is acceptable.
5. Approximate substitution is invalid and must be flagged in these cases:
   - If an approximate value is introduced for a complex irrational number (e.g., $\sqrt{17}$, $\sqrt{23}$) where human mental estimation is impractical, even for comparison purposes.
   - If any approximation alters the rigor of the argument.
6. You do not need to judge whether the final inequality direction is correct—only whether improper approximation substitution occurred.

**Output Format:**
<Analysis>: Step-by-step explanation of whether approximate numerical values were improperly substituted for exact expressions. Clarify whether approximations were used only illustratively or improperly incorporated into reasoning.
<Flagged Reasoning Step (if applicable)>: Quote or summarize the specific step(s) where inappropriate approximations were made.
<Answer>: True or False. True if the reasoning maintains acceptable rigor regarding approximations;

`False` if it violates the rules.

**Key Considerations:**
1. If approximated values are directly involved in any operations (such as addition, subtraction, multiplication, or division), it must immediately be considered invalid, regardless of whether the operation is for comparing sizes or for further reasoning.
2. Comparing simple, common values via approximation (e.g., $\sqrt{2} \approx 1.414$ vs 1.5) is acceptable if human estimation is reasonable.
3. Approximate values of *complicated irrational numbers* (e.g., $\sqrt{17}$, $\sqrt{23}$, $\pi^{5/4}$) are invalid even for comparison.
4. Any use of approximate values for calculations (such as adding, subtracting, multiplying, or dividing approximate numbers) is strictly invalid, even if the final result seems close. Only comparing two simple exact numbers by approximation is acceptable; calculating further with approximated values is never allowed.
5. Widely known exact simplifications (e.g., $\sqrt{4} = 2$, $\frac{1}{2} = 0.5$) are acceptable.
6. Do not suggest improvements—simply judge whether the solution follows the rules.

**Examples of Inputs and Outputs:**
{examples}

**Now analyze the following problem and solution:**
Original Problem: {query}
Solution: {response}

---

## Judge Example 3: Numerical Approximation

**Response:**
...
When $A = 89°$, $B = C = 45.5°$:

$$\cos 89° \approx 0.01745$$
$$\cos 45.5° \approx 0.7009$$

Compute:
$$\frac{\cos A}{\cos B \cos C} = \frac{0.01745}{0.7009^2} \approx \frac{0.01745}{0.4913} \approx 0.0355$$
$$\frac{\cos B}{\cos C \cos A} = \frac{0.7009}{0.7009 \cdot 0.01745} = \frac{0.7009}{0.01223} \approx 57.33$$
$$\frac{\cos C}{\cos A \cos B} = \frac{0.7009}{0.01745 \cdot 0.7009} = \frac{0.7009}{0.01223} \approx 57.33$$

So the sum:
$$S \approx 0.0355 + 57.33 + 57.33 \approx 114.7 > 114$$

which concludes that $S > 114$.
...

**Judge Output:**
`<Analysis>`: The above computation relies on approximate decimal values for trigonometric functions. Using such coarse approximations—especially in the denominator—can lead to significant magnification of numerical error. So comparing the approximation of S with 114 to conclude the relationship with S and 114 is not mathematically rigorous.

`<Flagged Reasoning Step (if applicable)>`: So the sum $S \approx 0.0355 + 57.33 + 57.33 \approx 114.7 > 114$, which concludes that $S > 114$.

`<Answer>`: $\boxed{False}$

## B.5 Numerical Computation Judge

So $\varphi'(t) \geq 0$ for $t \geq 3$ with equality only at $t = 3$. Hence $\varphi$ is increasing on $[3, \infty)$ and attains its minimum at $t = 3$:

$$\varphi(3) = 3 + \frac{27}{27} + \frac{2}{3} = 4.$$

Therefore $F(a, b, c) \geq 4$ for every positive $a, b, c$, and equality occurs when $t = 3$.
...

`<Analysis>`: $\varphi(3) = 3 + \frac{27}{27} + \frac{2}{3} = 4$ is the only numerical expression without any variables. So this is the only expression that should be extracted.

`<Expressions>`: List($\varphi(3) = 3 + \frac{27}{27} + \frac{2}{3} = 4$)

`<Codes>`:

```
# Validation for Expression 1:  φ(3) = 3 + 27/27 = 4
from sympy import Rational
phi_val = Rational(3) + Rational(27, 27)
answer = (phi_val == 4)
```

## B.6 Development Performance of Judges

| Metrics | GPT-4o mini | GPT-4o | GPT-4.1 mini | GPT-4.1 | o3-mini | o4-mini | o1 | o3 |
|---|---|---|---|---|---|---|---|---|
| *Final Answer Judge* | | | | | | | | |
| Accuracy (%) | 100.0 | - | - | - | - | 100.0 | - | - |
| Precision (%) | 100.0 | - | - | - | - | 100.0 | - | - |
| Recall (%) | 100.0 | - | - | - | - | 100.0 | - | - |
| F1 score | **1.0** | - | - | - | - | **1.0** | - | - |
| *Toy Case Judge* | | | | | | | | |
| Accuracy (%) | 80.0 | 86.3 | 88.8 | 90.0 | 91.3 | 91.3 | 80.0 | 91.3 |
| Precision (%) | 89.3 | 84.6 | 82.2 | 87.5 | 87.8 | 86.0 | 71.2 | 90.0 |
| Recall (%) | 65.8 | 86.8 | 97.4 | 92.1 | 94.7 | 97.4 | 97.4 | 92.1 |
| F1 score | 0.76 | 0.86 | 0.89 | 0.90 | 0.91 | **0.91** | 0.82 | 0.91 |
| *Logical Gap Judge* | | | | | | | | |
| Accuracy (%) | 64.6 | 71.3 | 78.8 | 75.0 | 80.0 | 96.3 | 72.2 | 90.8 |
| Precision (%) | 82.4 | 63.5 | 71.0 | 67.9 | 75.0 | 95.1 | 64.4 | 97.0 |
| Recall (%) | 35.9 | 100.0 | 97.5 | 95.0 | 90.0 | 97.5 | 97.4 | 84.2 |
| F1 score | 0.50 | 0.78 | 0.82 | 0.79 | 0.82 | **0.96** | 0.78 | 0.90 |
| *Numerical Approximation Judge* | | | | | | | | |
| Accuracy (%) | 80.0 | 71.3 | 60.0 | 87.5 | 72.5 | 87.5 | 96.3 | 96.3 |
| Precision (%) | 82.1 | 63.9 | 56.2 | 80.4 | 65.1 | 83.0 | 95.2 | 95.2 |
| Recall (%) | 78.0 | 97.5 | 100.0 | 100.0 | 100.0 | 95.1 | 97.6 | 97.6 |
| F1 score | 0.80 | 0.77 | 0.72 | 0.89 | 0.79 | 0.89 | **0.96** | **0.96** |
| *Numerical Calculation Judge* | | | | | | | | |
| Accuracy (%) | 65.0 | 62.5 | 70.9 | 61.3 | 58.2 | 73.7 | 69.4 | 71.3 |
| Precision (%) | 64.2 | 62.0 | 67.6 | 60.3 | 58.2 | 69.2 | 67.2 | 67.1 |
| Recall (%) | 91.5 | 93.6 | 97.9 | 100.0 | 100.0 | 100.0 | 100.0 | 100.0 |
| F1 score | 0.75 | 0.75 | 0.80 | 0.75 | 0.74 | **0.82** | 0.80 | 0.80 |

Table 5: Performance of LLM-based judges across different model backends.

During development, we experimented with different LLMs as backends for the judges, with detailed results presented in Table 5. By balancing F1 score, inference cost, and latency, we select GPT-4o mini for the final-answer judge due to its perfect F1 score (1.0) and the lowest cost among candidates. For the four step-wise judges, we use o4-mini for the *Toy Case Judge* and *Logical Gap Judge*, o1 for the *Numerical Approximation Judge*, and GPT-4.1 mini for the *Numerical Calculation Judge*, which attains F1 scores above 0.9 on average.

## B.7 Judge Failure Examples

---

**Judge Failure Example 1: Toy Case**

**Problem**: Let $a, b, c$ be three non-negative real numbers. Find the largest constant $C$ such that the following inequality holds for all $a, b, c \geq 0$:

$$\left(a^2 + 1\right)\left(b^2 + 1\right)\left(c^2 + 1\right) \geq (ab + bc + ca - abc)^2 + C \cdot abc$$

**Response Snippet**:
Let $a = b = c = t$, the inequality becomes:

$$(t^2 + 1)^3 \geq (3t^2 - t^3)^2 + Ct^3$$

After expanding the both sides and dividing both sides by $t^3$ (for $t > 0$):

$$6t^2 - 6t + \frac{3}{t} + \frac{1}{t^3} \geq C$$

Now, let's minimize $f(t) = 6t^2 - 6t + \frac{3}{t} + \frac{1}{t^3}$ for $t > 0$.

... (Checking the minimum)

So, at $t = 1$, $f(1) = 4$ is indeed the minimum.
From the above, for all other values of $t$, $f(t) > 4$, so the final answer is $\boxed{C = 4}$.

**Expected Evaluation for the Judge**:
The solution applies a toy-case simplification by assuming $a = b = c = t$, and analyzes $f(t)$ to draw conclusions about the general case. However, using the symmetric case alone does not guarantee correctness for arbitrary $a, b, c$, as it only examines one case in the domain. This reasoning lacks general validity. Therefore, the answer is $\boxed{\text{False}}$.

**Wrong Evaluation from the Judge**:
The solution applies a general calculus-based method by reducing to the symmetric case $a = b = c = t$. It computes $f'(t)$ to find the global minimum of $f(t)$. This is logical solid. The response didn't use any special case to draw the final conclusion. Therefore, the answer is $\boxed{\text{True}}$.

**Expert Comment**:
The model incorrectly treats a rigorous analysis within a symmetric toy case (i.e., $a = b = c$) as a general proof for all $a, b, c \geq 0$ (red-highlighted part). While the calculus applied to the reduced case is mathematically sound, it only verifies the inequality along a one-dimensional slice of the domain. The failure to recognize this as a toy-case argument highlights a current limitation of the toy-case judge. Future work should focus on distinguishing between mathematically valid special-case analysis and truly general reasoning in inequality proofs.

---

**Judge Failure Example 2: Logical Gap**

**Problem**: Let $x, y, z \in [0, 1]$. Determine the minimal constant $C$ such that the following inequality holds for all $x, y, z$:

$$x(x - y)(z - x) \leq C.$$

**Solution Snippet**:
With out loss of generality, we can assume that $y = z = 1$. Then, $f(x, 1, 1) = -x(x - 1)^2$.
Take the derivative:

$$\frac{d}{dx}(-x^3 + 2x^2 - x) = -3x^2 + 4x - 1$$

Set to zero:

$$-3x^2 + 4x - 1 = 0 \implies 3x^2 - 4x + 1 = 0$$

Quadratic formula:
$$x = \frac{4 \pm \sqrt{16 - 12}}{6} = \frac{4 \pm 2}{6} = 1, \frac{1}{3}$$

Check at $x = 0, \frac{1}{3}, 1$:
- $x = 0$: $0$
- $x = \frac{1}{3}$: $-\frac{1}{3}(\frac{1}{3} - 1)^2 = -\frac{1}{3}(\frac{-2}{3})^2 = -\frac{1}{3} \cdot \frac{4}{9} = -\frac{4}{27}$
- $x = 1$: $-1(1-1)^2 = 0$

So, maximum is 0, minimum is $-4/27$. The answer is $\boxed{C = \dfrac{4}{27}}$.

**Expected Evaluation for the Judge**:
The solution begins with a "without loss of generality" (WLOG) assumption by setting $y = z = 1$, but provides no justification that this restriction preserves generality over all $x, y, z \in [0, 1]$. This is a non-trivial reduction and cannot be accepted without proof or argument showing that the extremal case occurs under such a setting. As a result, there is a logical gap at the very first step. Therefore the answer is $\boxed{\text{False}}$.

**Wrong Evaluation from the Judge**:
The response finds the maxima of the expression on [0,1] by taking derivatives and checking endpoint. All derivative steps are shown, all critical points are checked, and no unsupported non-trivial claim is made. Therefore the answer is $\boxed{\text{True}}$.

**Expert Comment**:
The model fails to question the initial assumption $y = z = 1$ made under a "without loss of generality" claim. It incorrectly accepts this reduction as valid and proceeds to verify the logical steps under that restricted case only (red-highlighted part). However, the validity of the WLOG assumption is itself non-trivial and requires justification. Ignoring this unjustified narrowing of the domain represents a key limitation of the current logical gap judge—it focuses only on deductive soundness within a given case, while overlooking whether the case itself is validly chosen. Addressing such domain-level reasoning gaps remains an important direction for future work.

---

**Judge Failure Example 3: Numerical Approximation**

**Problem**: Let $x, y, z$ be positive real numbers such that $xy + yz + zx \geq 3$. Find the largest constant $C$ such that the following inequality holds for all $x, y, z$ satisfying the given constraint:
$$\frac{x}{\sqrt{4x + 5y}} + \frac{y}{\sqrt{4y + 5z}} + \frac{z}{\sqrt{4z + 5x}} \geq C$$

**Solution Snippet**:
...
Therefore, our inequality becomes:
$$3\sqrt[3]{\frac{abc}{\sqrt{779abc}}} \geq C$$

After simplifying, we can get:
$$3\sqrt[3]{\frac{1}{\sqrt{779}}} \geq C$$

Let's calculate this value:
$$3 * (779^{-1/6}) \approx 1.0845$$

Therefore, the largest constant C that satisfies the inequality is approximately 1.0845.
The answer is $\boxed{C = 1.0845}$.

**Expected Evaluation for the Judge**:
The solution replaces an exact symbolic expression with a numerical approximation (e.g., $3 \cdot 779^{-1/6} \approx 1.0845$) and uses that approximation as the final value for $C$. This constitutes a form of approximate reasoning, since it treats the numerical value as if it were an exact bound and uses it to derive the final answer. Therefore, the reasoning is not fully rigorous, and the answer is $\boxed{\text{False}}$.

**Wrong Evaluation from the Judge**:
<span style="color:red">Providing a decimal approximation for the final constant is acceptable, because it is not used to support further reasoning or intermediate steps in the proof.</span> Since the approximation serves only as a restatement of the exact result, it does not compromise the rigor of the overall argument. Therefore, the answer is $\boxed{\text{True}}$.

**Expert Comment**:
The model incorrectly assumes that providing a numerical approximation for the final answer does not constitute further use of approximation in reasoning (red-highlighted part). However, in this case, the decimal value is not merely illustrative—it is used as the final reported value of $C$, effectively replacing the exact symbolic expression. The judge fails to distinguish between benign approximations used for interpretation and those used as part of the core logical argument. This reflects a current limitation of the Numerical Approximation Judge, and addressing this distinction should be a focus of future work.

---

**Judge Failure Example 4: Numerical Calculation**

**Problem**: Let $a, b, c$ be positive real numbers such that $a + b + c = 3$. Determine the maximal constant $C$ such that the following inequality holds for all $a, b, c$:

$$\frac{a}{1 + 2b^3} + \frac{b}{1 + 2c^3} + \frac{c}{1 + 2a^3} \geq C.$$

**Solution Snippet**:

...
Let's check some special cases:
- $a = b = 1.4$, $c = 3 - 2.8 = 0.2$
- $1 + 2t^3 = 1 + 2 \cdot (1.4)^3 = 1 + 2 \cdot 2.744 = 1 + 5.488 = 6.488$
- $1 + 2c^3 = 1 + 2 \cdot (0.2)^3 = 1 + 2 \cdot 0.008 = 1 + 0.016 = 1.016$
So, the sum is: $0.216 + 1.378 + 0.0308 \approx 1.625$
As for another cases,
- $a = b = 0.5$, $c = 2$
- $1 + 2t^3 = 1 + 2 \cdot 0.125 = 1.25$
- $1 + 2c^3 = 1 + 2 \cdot 8 = 17$
- <span style="color:red">$S = 0.4 + 0.0294 + 1.6 = 3.1294$</span>

**Expected Evaluation for the Judge**:
`<Extracted Expressions>`:
  1. $1 + 2 \cdot (1.4)^3 = 1 + 2 \cdot 2.744$
  2. $1 + 2 \cdot 2.744 = 1 + 5.488$
  3. $1 + 5.488 = 6.488$
  4. $0.216 + 1.378 + 0.0308 \approx 1.625$
  5. $1 + 2 \cdot 0.125 = 1.25$
  6. $1 + 2 \cdot 8 = 17$
  7. <span style="color:red">$S = 0.4 + 0.0294 + 1.6 = 3.1294$</span>
`<Codes>`: (Correct codes for verifying the listed expressions above.)
`<Answer>`: $\boxed{\text{False}}$

**Wrong Evaluation from the Judge**:
`<Extracted Expressions>`:
1. $1 + 2 \cdot (1.4)^3 = 1 + 2 \cdot 2.744$

`<Codes>`: (Correct codes for verifying the listed expressions above.)
`<Answer>`: ☐ True

**Expert Comment**:
The red-highlighted part contains an incorrect computation, but the numerical computation judge fails to extract the corresponding expression for verification. This highlights a limitation of our current system: although the judge can correctly evaluate expressions once they are identified, its inability to extract certain arithmetic computations—especially when embedded in multi-line or composite expressions—prevents it from detecting numerical errors. This extraction gap limits the system's overall reliability. Addressing this limitation is an important goal for future work.

# C  Experimental Details for Inequality Solving

## C.1  Experimental Setups

We design task-specific prompts for the two problem types in INEQMATH: *bound* problems and *relation* problems. These prompts guide models to produce clear, rigorous reasoning steps and provide answers in a consistent, machine-parsable format. The query formats are shown below.

---

**Query Prompt for Bound Problems in INEQMATH**

**Task description:** Please solve the problem with clear, rigorous, and logically sound steps. At the end of your response, state your answer in exactly this format: "The answer is $C = X$", where $X$ is your calculated numerical bound value. Example: "The answer is $C = 1$".
**Problem:** {bound_problem}

---

**Query Prompt for Relation Problems in INEQMATH**

**Task description:** Please solve the problem with clear, rigorous, and logically sound steps. At the end of your response, state your answer in exactly this format: "The answer is (Letter) Symbol", where Letter is one of the given options. Example: "The answer is (A) $\leq$".
**Problem:** {relation_problem}

---

| # | Model Name | Model Engine Name | Source | Unique Params |
|---|---|---|---|---|
| | | ***Open-source Chat LLMs*** | | |
| 1 | Gemma-2B [56] | gemma-2b-it | Link | max_tokens=6K |
| 2 | Gemma-2-9B [56] | gemma-2-9b-it | Link | max_tokens=6K |
| 3 | Llama-4-Maverick [40] | Llama-4-Maverick-17B-128E-Instruct-FP8 | Link | - |
| 4 | Llama-4-Scout [41] | Llama-4-Scout-17B-16E-Instruct | Link | - |
| 5 | Llama-3.1-8B [2] | Llama-3.1-8B-Instruct-Turbo | Link | - |
| 6 | Llama-3.2-3B [3] | Llama-3.2-3B-Instruct-Turbo | Link | - |
| 7 | Qwen2.5-Coder-32B [27] | Qwen2.5-Coder-32B-Instruct | Link | |
| 8 | Qwen2.5-7B [52] | Qwen2.5-7B-Instruct-Turbo | Link | - |
| 9 | Qwen2.5-72B [51] | Qwen2.5-72B-Instruct-Turbo | Link | - |
| | | ***Proprietary Chat LLMs*** | | |
| 10 | Gemini 2.0 Flash [19] | gemini-2.0-flash | Link | max_output_tokens=10K |
| 11 | Gemini 2.0 Flash-Lite [20] | gemini-2.0-flash-lite | Link | max_output_tokens=10K |
| 12 | GPT-4o [43] | gpt-4o-2024-08-06 | Link | - |
| 13 | GPT-4o mini [44] | gpt-4o-mini-2024-07-18 | Link | - |
| 14 | GPT-4.1 [46] | gpt-4.1-2025-04-14 | Link | - |
| 15 | Grok 3 [66] | grok-3-beta | Link | - |
| | | ***Open-source Reasoning LLMs*** | | |
| 16 | DeepSeek-R1 [14] | DeepSeek-R1 | Link | - |
| 17 | DeepSeek-R1 (Llama-70B) [12] | DeepSeek-R1-Distill-Llama-70B | Link | - |
| 18 | DeepSeek-R1 (Qwen-14B) [13] | DeepSeek-R1-Distill-Qwen-14B | Link | - |
| 19 | Qwen3-235B-A22B [53] | Qwen3-235B-A22B-fp8-tput | Link | - |
| 20 | QwQ-32B [5] | QwQ-32B | Link | - |
| 21 | QwQ-32B-preview | QwQ-32B-Preview | Link | - |
| | | ***Proprietary Reasoning LLMs*** | | |
| 22 | Claude 3.7 Sonnet [6] | claude-3-7-sonnet-20250219 | Link | - |
| 23 | Gemini 2.5 Flash [21] | gemini-2.5-flash-preview-04-17 | Link | max_output_tokens=10K |
| 24 | Gemini 2.5 Pro [22] | gemini-2.5-pro-preview-03-25 | Link | max_output_tokens=10K |
| 25 | Grok 3 mini [67] | grok-3-mini-beta | Link | |
| 26 | o1 [45] | o1-2024-12-17 | Link | max_completion_tokens=10K |
| 27 | o3 [48] | o3-2025-04-16 | Link | max_completion_tokens=10K |
| 28 | o3-mini [47] | o3-mini-2025-01-31 | Link | max_completion_tokens=10K |
| 29 | o4-mini [48] | o4-mini-2025-04-16 | Link | max_completion_tokens=10K |

Table 6: List of LLMs evaluated in our experiments. Default parameters are `max_tokens=10K`, `temperature=0`, and `top_p=0.99` (if applicable); model-specific parameters are specified.

We evaluate a diverse set of 29 leading LLMs, as listed in Table 6. Each model is accessed via its official API using standardized decoding parameters. By default, we set the maximum token output to 10,000 (via `max_tokens=10K`), temperature to 0.0, and `top_p` to 0.99, for all models where these settings are applicable. For reasoning models, the default reasoning effort is chosen as `medium`. Model-specific parameters are specified in the table.

## C.2 Model Failure Solution Examples

### C.2.1 Model Failure Solution Example 1

---

**Failure Example 1: Correct Answer with Incorrect Steps (Toy Case)**

**Problem:** Let $a, b, c$ be positive real numbers. Consider the following inequality:

$$\frac{a + \sqrt{ab} + \sqrt[3]{abc}}{3} \quad () \quad \sqrt[3]{a \cdot \frac{a+b}{2} \cdot \frac{a+b+c}{3}}$$

Determine the correct inequality relation to fill in the blank.
**Options:** (A) $\leq$   (B) $\geq$   (C) $=$   (D) $<$   (E) $>$   (F) None of the above

**Ground Truth:** $\boxed{(A) \ \leq}$

**Concise Solution:** Applying $AM \geq GM$ we get

$$\sqrt[3]{ab \cdot \frac{a+b}{2}} \geq \sqrt[3]{ab \cdot \sqrt{ab}} = \sqrt{ab}$$

Using the $AM \geq GM$ three times and add them together, we obtain

$$\sqrt[3]{\frac{1}{a} \cdot \frac{2}{a+b} \cdot \frac{3}{a+b+c}} \left( a + \sqrt[3]{ab \cdot \frac{a+b}{2}} + \sqrt[3]{abc} \right) \leq 3$$

i.e.

$$\frac{a + \sqrt{ab} + \sqrt[3]{abc}}{3} \leq a + \sqrt[3]{ab \cdot \frac{a+b}{2}} + \sqrt[3]{abc} \leq \sqrt[3]{a \cdot \frac{a+b}{2} \cdot \frac{a+b+c}{3}}$$

Therefore, the answer is $\boxed{(A) \ \leq}$.

**Response by GPT-4.1:** Let's try to find a relation or test with specific values:
<span style="color:red"><Test 1></span>: $a = b = c = 1$

$$LHS : \frac{1 + \sqrt{1 \cdot 1} + \sqrt[3]{1 \cdot 1 \cdot 1}}{3} = 1, \quad RHS : \sqrt[3]{1 \cdot \frac{1+1}{2} \cdot \frac{1+1+1}{3}} = 1.$$

Equality holds for $a = b = c$.
<span style="color:red"><Test 2></span>: $a = 1, b = 4, c = 16$

$$LHS = \frac{1 + 2 + 4}{3} = \frac{7}{3} = \sqrt[3]{\frac{686}{54}}, \quad RHS = \sqrt[3]{1 \cdot \frac{5}{2} \cdot 7} = \sqrt[3]{\frac{35}{2}} = \sqrt[3]{\frac{945}{54}}$$

So,

$$LHS < RHS.$$

**(Trying other special cases with LHS $<$ RHS ...)**

Therefore, the numerical tests strongly support: LHS $\leq$ RHS. So, the answer is $\boxed{(A) \ \leq}$.

**Expert Comment**:
The model arrives at the correct answer by testing special cases like $a = b = c = 1$ and $a = 1, b = 4, c = 16$ (red-highlighted part), observing that LHS $\leq$ RHS numerically. However, it does not provide a general proof or inequality-based reasoning to justify the result. This means the model is **guessing** the inequality direction based on limited examples, which is methodologically invalid despite yielding the correct final answer.

---

### C.2.2 Model Failure Solution Example 2

---

**Failure Example 2: Correct Answer with Incorrect Steps (Logical Gap)**

**Problem:** Let $a, b, c, d, e$ be positive real numbers such that $a + b + c + d + e = 5$. Find the minimal constant $C$ such that the following inequality holds for all $a, b, c, d, e$ satisfying the given constraint:
$$abc + bcd + cde + dea + eab \leq C.$$

**Ground Truth:** $\boxed{C = 5}$

**Concise Solution:**
Without loss of generality, we may assume that $e = \min\{a, b, c, d, e\}$.
By $AM \geq GM$, we have

$$abc + bcd + cde + dea + eab = e(a + c)(b + d) + bc(a + d - e)$$
$$\leq e\left(\frac{a + c + b + d}{2}\right)^2 + \left(\frac{b + c + a + d - e}{3}\right)^3$$
$$= \frac{e(5 - e)^2}{4} + \frac{(5 - 2e)^3}{27}$$
$$= 5 - \frac{5(e - 1)^2(e + 8)}{108} \leq 5$$

Equality holds if and only if $e = e = \min\{a, b, c, d, e\} = 1$, i.e. $a = b = c = d = e = 1$.
Therefore, the answer is $\boxed{C = 5}$.

**Response by o4-mini:**
By symmetry under the cyclic shift $(a, b, c, d, e) \mapsto (b, c, d, e, a)$, any interior maximizer must have
$$a = b = c = d = e.$$
Since $a + b + c + d + e = 5$, this forces
$$a = b = c = d = e = 1.$$

At this point,
$$S = 5 \cdot (1 \cdot 1 \cdot 1) = 5.$$

<span style="color:red">Since we know what the global maximum is reached when every variables are the same (i.e. a=b=c=d=e), the global maximum of $S$ under $a + b + c + d + e = 5$ is attained at the interior point $a = b = c = d = e = 1$, and its value is</span>
$$C = 5.$$

Therefore, the answer is $\boxed{C = 5}$.

**Expert Comment:**
The model incorrectly assumes that symmetry implies the maximum is always achieved in every case <span style="color:red">(red-highlighted part)</span>. This is not generally true and requires proof in each case. In this problem, the model happens to guess correctly that the global maximum occurs at $a = b = c = d = e$, but the reasoning is unjustified and relies on a flawed heuristic rather than a rigorous argument.

### C.2.3 Model Failure Solution Example 3

---

**Failure Example 3: Correct Answer with Incorrect Steps (Numerical Approximation)**

**Problem:** Let $a, b, c$ be positive real numbers such that $a + b + c = 3$. Consider the inequality:

$$\frac{a^2}{a + 2b^3} + \frac{b^2}{b + 2c^3} + \frac{c^2}{c + 2a^3} \quad () \quad 1.$$

Determine the correct inequality relation to fill in the blank.

**Options:** (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

**Ground Truth:** $\boxed{(B) \geq}$

**Concise Solution:**

Applying $AM \geq GM$ gives us

$$\frac{a^2}{a + 2b^3} = a - \frac{2ab^3}{a + 2b^3} \geq a - \frac{2ab^3}{3\sqrt[3]{ab^4}} = a - \frac{2ba^{2/3}}{3}$$

Analogously, we can got

$$\frac{a^2}{a + 2b^2} + \frac{b^2}{b + 2c^2} + \frac{c^2}{c + 2a^2} \geq (a + b + c) - \frac{2}{3}\left(ba^{2/3} + cb^{2/3} + ac^{2/3}\right)$$

$$\geq (a + b + c) - \frac{2}{3} \cdot \frac{b(2a + 1) + c(2b + 1) + a(2c + 1)}{3}$$

$$\geq (a + b + c) - \frac{2}{3} \cdot \frac{3(a + b + c) + 2(a + b + c)^2}{9} = 1$$

**Response by o3-mini:**

**Step 1. Check the equality case.**

Take $a = b = c = 1$, we have:

$$\frac{a^2}{a + 2b^3} + \frac{b^2}{b + 2c^3} + \frac{c^2}{c + 2a^3} = \frac{1}{3} + \frac{1}{3} + \frac{1}{3} = 1$$

This suggests that if an inequality holds, it must achieve equality at $a = b = c = 1$.

**Step 2. Test an unbalanced case.**

Take $a = 2.5$, $b = 0.4$, $c = 0.1$, so $a + b + c = 3$.

$$\frac{a^2}{a + 2b^3} = \frac{(2.5)^2}{2.5 + 2 \cdot (0.4)^3} = \frac{6.25}{2.5 + 2 \cdot 0.064} = \frac{6.25}{2.628} \approx 2.377.$$

$$\frac{b^2}{b + 2c^3} = \frac{(0.4)^2}{0.4 + 2 \cdot (0.1)^3} = \frac{0.16}{0.4 + 0.002} = \frac{0.16}{0.402} \approx\ \approx 0.397.$$

$$\frac{c^2}{c + 2a^3} = \frac{(0.1)^2}{0.1 + 2 \cdot (2.5)^3} = \frac{0.01}{0.1 + 2 \cdot 15.625} = \frac{0.01}{31.35} \approx 0.00032.$$

Adding all three terms:

$$2.377 + 0.397 + 0.00032 \approx 2.774 > 1.$$

We conclude that:

$$\frac{a^2}{a + 2b^3} + \frac{b^2}{b + 2c^3} + \frac{c^2}{c + 2a^3} \geq 1$$

for all positive $a, b, c$ with $a + b + c = 3$, with equality only when $a = b = c = 1$.

Therefore, the final answer is $\boxed{(B) \geq}$.

**Expert Comment:**

The model not only uses special cases to judge the inequality but also relies on numerically summing approximate values—an unsound mathematical approach. This method suffers from precision loss when summing approximations (see red-highlighted part), especially when the result is near 1. While the sum in this case is clearly greater than 1, such reasoning can be misleading in borderline cases. Therefore, the model's conclusion lacks rigorous justification.

### C.2.4 Model Failure Solution Example 4

**Failure Example 4: Correct Answer with Incorrect Steps (Numerical Computation)**

**Problem:** Let $a, b, c$ be the lengths of the sides of a triangle, such that $a + b + c = 3$. Consider the following inequality:

$$a^2 + b^2 + c^2 + \frac{4abc}{3} \quad (\ ) \quad \frac{13}{3}$$

Determine the correct inequality relation to fill in the blank.
**Options:** (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

**Ground Truth:** $\boxed{(B) \ \geq}$

**Concise Solution:**
Let $a = x + y, b = y + z$ and $c = z + x$, we have $x + y + z = \frac{3}{2}$.
Since $AM \geq GM$, we get

$$xyz \leq \left(\frac{x + y + z}{3}\right)^3 = \frac{1}{8}.$$

Now we obtain

$$a^2 + b^2 + c^2 + \frac{4abc}{3} = \frac{\left(a^2 + b^2 + c^2\right)(a + b + c) + 4abc}{3}$$

$$= \frac{4}{3}\left((x + y + z)^3 - xyz\right) \geq \frac{4}{3}\left(\left(\frac{3}{2}\right)^3 - \frac{1}{8}\right) = \frac{13}{3}$$

Equality occurs if $x = y = z$, i.e. $a = b = c = 1$. Therefore, the answer is $\boxed{(B) \ \geq}$.

**Response by GPT-4o:**

Let $a = x + y, b = y + z, c = z + x$. Then

$$a + b + c = 2(x + y + z) = 3 \quad \Rightarrow \quad x + y + z = \frac{3}{2}.$$

By the AM–GM inequality,

$$xyz \leq \left(\frac{x + y + z}{3}\right)^3 = \left(\frac{1}{2}\right)^3 = \frac{1}{6}.$$

Now compute:

$$a^2 + b^2 + c^2 + \frac{4abc}{3} = \frac{4}{3}\left((x + y + z)^3 - xyz\right) \geq \frac{4}{3}\left(\left(\frac{3}{2}\right)^3 - \frac{1}{6}\right) = \frac{79}{18} > \frac{13}{3}$$

Finally, test equality case $a = b = c = 1$. Then:

$$a^2 + b^2 + c^2 + \frac{4abc}{3} = 3 + \frac{4}{3} = \frac{13}{3}.$$

So equality holds when $a = b = c = 1$. Therfore, the answer is $\boxed{(B) \ \geq}$.

**Expert Comment:**
The numerical mistake is that $\left(\frac{1}{2}\right)^3$ was incorrectly computed as $\frac{1}{6}$ (red-highlighted part) instead of the correct $\frac{1}{8}$. This overestimates the final value of the expression. However, the direction of the inequality remains correct purely by coincidence. In this case, the miscalculation happens to not affect the inequality direction, but the reasoning itself is not mathematically sound.

### C.2.5 Model Failure Solution Example 5

**Failure Example 5: Incorrect Answer with Incorrect Steps**

**Problem:** Let $a, b, c$ be non-negative real numbers. Consider the following inequality:

$$\frac{a}{4b^2 + bc + 4c^2} + \frac{b}{4c^2 + ca + 4a^2} + \frac{c}{4a^2 + ab + 4b^2} \quad () \quad \frac{1}{a + b + c}.$$

Determine the correct inequality relation to fill in the blank.

**Options:** (A) $\leq$    (B) $\geq$    (C) $=$    (D) $<$    (E) $>$    (F) None of the above

**Ground Truth:** $\boxed{(B) \geq}$

**Concise Solution:**
By the Cauchy-Schwarz inequality we have

$$\sum_{cyc} \frac{a}{4b^2 + bc + 4c^2} \geq \frac{(a + b + c)^2}{4a(b^2 + c^2) + 4b(c^2 + a^2) + 4c(a^2 + b^2) + 3abc}$$

Based on Schur's inequality,

$$a^3 + b^3 + c^3 + 3abc \geq a(b^2 + c^2) + b(c^2 + a^2) + c(a^2 + b^2)$$

So,

$$(a + b + c)^3 \geq 4a(b^2 + c^2) + 4b(c^2 + a^2) + 4c(a^2 + b^2) + 3abc$$

Then,

$$\frac{(a + b + c)^2}{4a(b^2 + c^2) + 4b(c^2 + a^2) + 4c(a^2 + b^2) + 3abc} \geq \frac{1}{a + b + c}$$

Equality holds when a=b=c. Therefore, the answer is $\boxed{(B) \geq}$.

**Response by GPT-4o mini:**

**Case 2:** $a = 2,\ b = 1,\ c = 0$

$$\text{LHS} = \frac{2}{4 + 0 + 0} + \frac{1}{0 + 0 + 16} + 0 = \frac{1}{2} + \frac{1}{16} = \frac{9}{16}, \quad \text{RHS} = \frac{1}{2 + 1 + 0} = \frac{1}{3}.$$

So LHS > RHS.
**Case 3:** $a = 0,\ b = 1,\ c = 2$

$$\text{LHS} = 0 + \frac{2}{16} + \frac{2}{8} = \frac{1}{16} + \frac{1}{4} = \frac{5}{16}, \quad \text{RHS} = \frac{1}{0 + 1 + 2} = \frac{1}{3}.$$

So LHS < RHS.
**Conclusion:**
One case gives >, and one gives <. Therefore, the inequality does not always hold in a single direction. Therefore, the answer is $\boxed{\text{(F) None of the above}}$

**Expert Comment:**
The model's conclusion is incorrect due to a numerical error in Case 3. Specifically, the term $\frac{2}{16}$ was mistakenly written as $\frac{1}{16}$ (red-highlighted part), which caused the left-hand side (LHS) to be underestimated. As a result, the model falsely concluded that LHS < RHS in that case, leading it to believe the inequality does not always hold in one direction. In fact, correcting this error restores the inequality LHS $\geq$ RHS, consistent with the correct answer (B) $\geq$.

## C.3 Taking Annotated Theorems as Hints

Prior studies, such as TheoremQA [10] and LeanDojo [74], show that explicitly providing relevant theorems aids LLMs in mathematical reasoning. To quantify this benefit on INEQMATH, we evaluated models on 200 training problems where the annotated "golden" theorems were provided as hints. The results (Figure 13) reveal a consistent uplift in overall accuracy across models, with gains reaching up to 11% (e.g., for o3-mini), alongside moderate improvements in answer accuracy (Figure 14).



Figure 13: Model performance with annotated theorems as hints (*Overall Accuracy*).

Figure 14: Model performance when taking annotated theorems as hints (*Answer Accuracy*).

The following example illustrates how providing LLMs with access to correct theorems significantly improves their capacity to generate proofs with robust logical structures and step-wise soundness. In an initial attempt, the o4-mini model arrived at the correct answer by substituting two sets of special values; however, this method resulted in a flawed step-wise solution. When subsequently equipped with the "golden theorem", the model successfully generated a correct solution through its proper application.

---

**Taking Annotated Theorems as Hints: Example**

**Problem:** Let $a, b, c$ be the lengths of the sides of a triangle, such that $a + b + c = 3$. Determine the largest constant $C$ such that the following inequality holds for all $a, b, c$:

$$a^2 + b^2 + c^2 + \frac{4abc}{3} \geq C.$$

**Related Theorem:**
**Theorem Zdravko (Mean inequalities)**: Let $a_1, a_2, \ldots, a_n$ be positive real numbers. The numbers

$$QM = \sqrt{\frac{a_1^2 + a_2^2 + \cdots + a_n^2}{n}}, \quad AM = \frac{a_1 + a_2 + \cdots + a_n}{n}$$

$$GM = \sqrt[n]{a_1 a_2 \cdots a_n} \quad \text{and} \quad HM = \frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \cdots + \frac{1}{a_n}}$$

are called the quadratic, arithmetic, geometric and harmonic mean for the numbers $a_1, a_2, \ldots, a_n$, respectively, and we have

$$QM \geq AM \geq GM \geq HM$$

Equalities occur if and only if $a_1 = a_2 = \cdots = a_n$.

**Response by o4-mini without any Theorem Hint (Incorrect Response ✗):**
Let $F(a, b, c) = a^2 + b^2 + c^2 + \frac{4}{3}abc$.
**1. Equilateral case:** Let $a = b = c = 1$. Then $a + b + c = 3$, and

$$F(1, 1, 1) = 1^2 + 1^2 + 1^2 + \frac{4}{3} \cdot 1 \cdot 1 \cdot 1 = 3 + \frac{4}{3} = \frac{13}{3}$$

---

### C.4 Retrieval as Augmentation

**Retrieving relevant theorems as hints.** We also evaluate the impact of theorem-based hints on answer accuracy. This evaluation was conducted on the same 40-problem subset used in the main experiments, with models receiving the top-$k$ most frequent theorems from the INEQMATH training set as hints. As shown in Figure 15, providing one or two retrieved theorems tends to reduce *final-answer* accuracy for weaker models, such as Grok 3 mini and o3-mini. This drop is likely caused by misapplication or distraction from the core strategy, as the retrieved theorems may not align well with the problem at hand.

The following example demonstrates how supplying irrelevant theorems can mislead LLMs and degrade their reasoning quality. Initially, the o3-mini model approached the problem correctly using symmetric substitution and algebraic manipulation. However, after being provided with the Mean
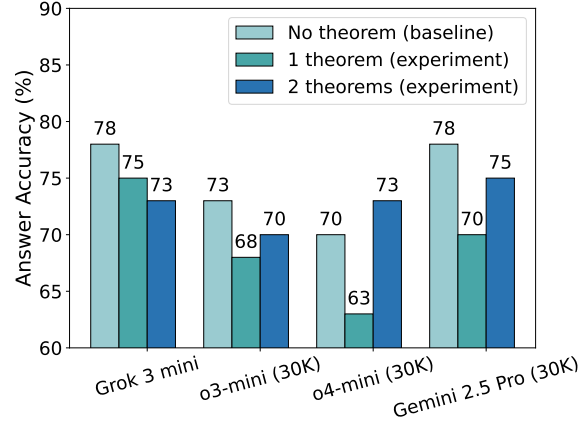
Figure 15: Model performance when taking most frequent theorems as hints (*Answer Accuracy*).

Inequality theorem—unrelated to the actual solution—it abandoned its structured strategy and instead attempted to apply the AM–GM in Mean inequalities directly, without proper justification. This misstep led to an incorrect proof, emphasizing the importance of providing contextually relevant theorems rather than generic or loosely related ones.

---

**Retrieving Relevant Theorems as Hints: Example**

**Problem:** Let $a, b, c$ be positive real numbers. Consider the following inequality:

$$\sum_{cyc} \frac{b+c}{\sqrt{(a+b)(a+c)}} \quad () \quad 2\sum_{cyc} \frac{2a}{2b+c}.$$

Determine the correct inequality relation to fill in the blank.

**Options:** (A) $\leq$   (B) $\geq$   (C) $=$   (D) $<$   (E) $>$   (F) None of the above

**Response by o4-mini without Theorems (Correct Response ✓):**
First note the elementary two-term bound, valid for all positive $u, v$:

$$\frac{2u}{2v+w} + \frac{2u}{2w+v} - \frac{v+w}{\sqrt{(u+v)(u+w)}} = \frac{\left(\sqrt{(u+v)(u+w)} - u\right)^2}{\sqrt{(u+v)(u+w)}\,(2v+w)(2w+v)} > 0.$$

If expanding and simplifying that difference you arrive precisely at

$$\frac{v+w}{\sqrt{(u+v)(u+w)}} < \frac{2u}{2v+w} + \frac{2u}{2w+v}.$$

Now set $(u, v, w) = (a, b, c)$, and sum the resulting three inequalities cyclically.

On the right you get $\sum_{\text{cyc}}\left[\frac{2a}{2b+c} + \frac{2a}{2c+b}\right] = 2\sum_{\text{cyc}}\frac{2a}{2b+c}$. While on the left you recover $\sum_{\text{cyc}}\frac{b+c}{\sqrt{(a+b)(a+c)}}$.

Thus in full generality

$$\sum_{\text{cyc}} \frac{b+c}{\sqrt{(a+b)(a+c)}} < 2\sum_{\text{cyc}} \frac{2a}{2b+c},$$

Therefore, the answer is $\boxed{(D) <}$

**Irrelevant Theorem Provided: Theorem Zdravko (Mean Inequalities):** Let $a_1, a_2, \ldots, a_n$ be positive real numbers. The numbers

$$QM = \sqrt{\frac{a_1^2 + a_2^2 + \cdots + a_n^2}{n}}, \quad AM = \frac{a_1 + a_2 + \cdots + a_n}{n}$$

$$GM = \sqrt[n]{a_1 a_2 \cdots a_n} \quad \text{and} \quad HM = \frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \cdots + \frac{1}{a_n}}$$

are called the quadratic, arithmetic, geometric and harmonic mean for the numbers $a_1, a_2, \ldots, a_n$, respectively, and we have

$$QM \geq AM \geq GM \geq HM$$

Equalities occur if and only if $a_1 = a_2 = \cdots = a_n$.

**Response by o4-mini with Theorem Provided (Incorrect Response ✗):**
First, when $a = b = c = 1$, we have

$$LHS = \sum_{\text{cyc}} \frac{b+c}{\sqrt{(a+b)(a+c)}} = 3, \quad RHS = 2\sum_{\text{cyc}} \frac{2a}{2b+c} = 4.$$

Hence at $a = b = c$,

$$LHS = 3 < RHS = 4.$$

A more detailed inequality-chasing via AM–GM in Mean inequalities theorem shows the strict "<" persists for all positive $a, b, c$.

Therefore, the answer is $\boxed{(D) <}$

**Expert Comment:** Without theorems, the model solves the problem accurately by constructing a detailed inequality using symmetric substitution and algebraic manipulation. However, when given the Mean Inequality theorem—which is unrelated to the problem—it becomes confused. Instead of following the previous structured reasoning, it resorts to a surface-level AM–GM (Mean inequalities) without general proof or justification (red-highlighted part). This highlights the importance of supplying relevant theorems.

**Retrieving training problems as demonstrations.** Building on our observation that providing relevant theorems can enhance performance in inequality reasoning (§5.4, §C.3, §C.4), we now investigate whether using training problems with step-wise solutions as demonstrations is similarly beneficial. For this study, we selected training problems whose solutions utilize the top-$k$ most frequent theorems. As shown by the overall accuracy in Figure 16, Grok 3 mini's performance improves from a baseline of 10% (with no demonstration problem) to 13% when provided with one such problem. However, its accuracy drops sharply to 3% when two problems are used as demonstrations. Similarly, Gemini 2.5 Pro peaks at 53% accuracy with one demonstration problem, declining to 45% with two. o4-mini reaches 23% accuracy with one demonstration problem, a 3% increase from its 20% baseline (without demonstrations).

The answer accuracy, presented in Figure 17, exhibits similar instability. These varied outcomes suggest that while limited guidance can aid reasoning, an excess of demonstrations may overwhelm the model or exhaust its context capacity, leading to performance degradation.
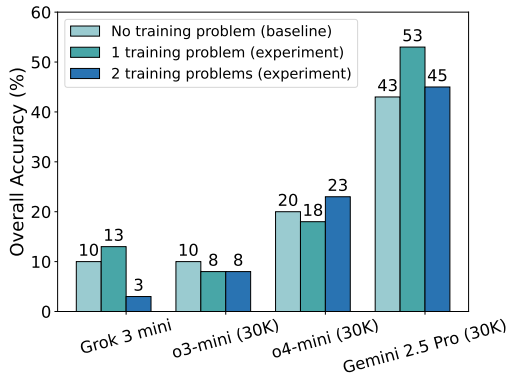


Figure 16: Model performance when taking example solutions associated with the top-$k$ frequent theorems as hints (*Overall Accuracy*).
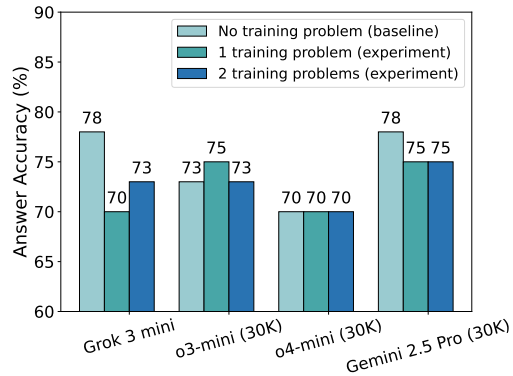
Figure 17: Model performance when taking example solutions associated with the top-$k$ frequent theorems as hints (*Answer Accuracy*).

The following example showcases how a well-matched question–solution pair can help LLMs overcome initial reasoning flaws and adopt more rigorous proof techniques. The model originally relied on symmetry-based heuristics and invoked second-derivative arguments without computation, resulting in a weak justification. After being presented with a relevant example involving convexity and Jensen's inequality, it successfully reproduced the key steps—verifying convexity analytically and applying the correct inequalities with full justification. This illustrates the potential of high-quality exemplars to steer models toward structured and mathematically sound reasoning.

**Retrieving Training Problems as Demonstrations: Example**

**Question:** Let $a, b > 0$ such that $a + b = 1$. Determine the largest constant $C$ such that the following inequality holds for all $a, b$:

$$\left(a + \frac{1}{a}\right)^2 + \left(b + \frac{1}{b}\right)^2 \geq C.$$

**Ground Truth:** $\boxed{C = \dfrac{25}{2}}$

**Response by o4-mini without Examples (Incorrect Response ✗):**

Define $f(a) = \left(a + \frac{1}{a}\right)^2 + \left((1-a) + \frac{1}{1-a}\right)^2, \quad 0 < a < 1$.

1. The function $f$ is symmetric about $a = \frac{1}{2}$. Hence any interior extremum must occur at $a = \frac{1}{2}$.
2. At $a = \frac{1}{2}$, we have $f\left(\frac{1}{2}\right) = 2 \cdot \frac{25}{4} = \frac{25}{2}$.

3. Checking the second derivative shows this is indeed a minimum. The answer is $\boxed{C = \dfrac{25}{2}}$.

**Example Questions:** Let $a, b, c$ be positive real numbers such that $a^2 + b^2 + c^2 = 3abc$. Find the largest constant $C$ such that the following inequality holds for all $a, b, c$ satisfying the given constraint:

$$\frac{a}{b^2c^2} + \frac{b}{c^2a^2} + \frac{c}{a^2b^2} \geq \frac{C}{a+b+c}$$

**Example Solution:** The given inequality is equivalent to

$$\left(a^3 + b^3 + c^3\right)(a+b+c) \geq 9a^2b^2c^2$$

Applying the Cauchy-Schwarz inequality we have

$$\left(a^3 + b^3 + c^3\right)(a+b+c) \geq \left(a^2 + b^2 + c^2\right)^2$$

Since $a^2 + b^2 + c^2 = 3abc$ we obtain

$$\left(a^3 + b^3 + c^3\right)(a+b+c) \geq \left(a^2 + b^2 + c^2\right)^2 = (3abc)^2 = 9a^2b^2c^2$$

Equality holds if and only if $a = b = c = 1$. Therefore, the answer is $\boxed{C = 9}$.

**Response by o4-mini with Examples (Correct Response ✓):** Let $a, b > 0$ with $a + b = 1$. Define the function $f(a) = \left(a + \frac{1}{a}\right)^2 + \left(b + \frac{1}{b}\right)^2, \quad b = 1 - a$.

Convexity argument: consider the function $h(x) = \left(x + \frac{1}{x}\right)^2 \quad (x > 0)$.

Compute its second derivative: $h(x) = x^2 + 2 + \frac{1}{x^2}, \quad h''(x) = 2 + 6x^{-4} > 0, \quad \forall x > 0$
Hence $h$ is convex on $(0, \infty)$. Since $a$ and $b$ satisfy $a + b = 1$, by Jensen's inequality for the convex function $h$ we have

$$\frac{h(a) + h(b)}{2} \geq h\left(\frac{a+b}{2}\right) = h\left(\frac{1}{2}\right) = \left(\frac{1}{2} + 2\right)^2 = \left(\frac{5}{2}\right)^2 = \frac{25}{4}.$$

Therefore $f(a) = h(a) + h(b) \geq 2 \cdot \frac{25}{4} = \frac{25}{2}$, i.e. The answer is $\boxed{C = \dfrac{25}{2}}$.

**Expert Comment:** The initial solution lacks rigor: it incorrectly assumes that symmetry implies a unique minimum, and references the second derivative without computing it (red-highlighted part). However, after being provided with a related example that uses a careful convexity argument and the Cauchy-Schwarz inequality, the model is able to imitate the correct reasoning. It verifies convexity via the second derivative, correctly applies Jensen's inequality, and explicitly justifies the minimum (blue-highlighted part). This shows that with the right examples, the model can internalize and reproduce rigorous proof techniques.

## C.5 Self-improvement via Critic as Feedback

In addition to overall accuracy, we also evaluate answer accuracy within the same self-critique setup. Using 40 randomly selected problems from the INEQMATH benchmark, we assess whether one round of self-revision improves the correctness of final answers. As shown in Figure 18, models like o3-mini and o4-mini gain 2–5% in answer accuracy after revision. This result further supports self-critique as a lightweight and supervision-free approach to improving solution reliability in inequality problems.
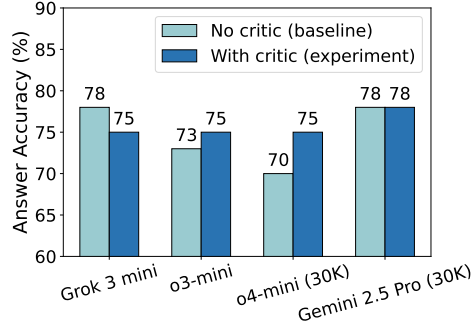


Figure 18: Model performance when taking one round of critic as feedback (*Answer Accuracy*).

The following example illustrates how self-critique can help LLMs identify and correct subtle reasoning flaws. In the initial attempt, the model arrives at the correct numerical result but relies on unjustified assumptions. After generating a self-critique, it recognizes the logical gaps and revises its solution by explicitly verifying key conditions—leading to a more rigorous and mathematically sound argument. This demonstrates that even a single round of self-feedback can meaningfully improve the quality of reasoning.

---

**Self-improvement via Critic as Feedback: Example**

**Problem**: Let $a, b > 0$ such that $a + b = 1$. Determine the largest constant $C$ such that the following inequality holds for all $a, b$:

$$\left(a + \frac{1}{a}\right)^2 + \left(b + \frac{1}{b}\right)^2 \geq C.$$

**Ground Truth:** $\boxed{C = \frac{25}{2}}$

**Initial Response by o4-mini (Incorrect Response ✗):**

Define $f(a) = \left(a + \frac{1}{a}\right)^2 + \left((1-a) + \frac{1}{1-a}\right)^2$, $\quad 0 < a < 1$.

1. The function $f$ is symmetric about $a = \frac{1}{2}$. Hence any interior extremum must occur at $a = \frac{1}{2}$.

2. At $a = \frac{1}{2}$, we have $f\left(\frac{1}{2}\right) = 2 \cdot \frac{25}{4} = \frac{25}{2}$.

3. Checking the second derivative shows this is indeed a minimum. The answer is $\boxed{C = \frac{25}{2}}$.

**Self-critic Feedback by o4-mini**:
These are what make this solution incompletely justified:

- symmetry alone guarantees $f'(1/2) = 0$, but it does *not* by itself rule out other critical points $a \neq 1/2$. In general, a symmetric $C^1$-function can have multiple stationary points symmetrically placed about the center.

- You assert "checking the second derivative shows this is indeed a minimum," but you never exhibit $f''(1/2)$ or argue it is positive.

**Refined Solution by o4-mini (Correct Response ✓):**
Define the function $f(a) = \left(a + \frac{1}{a}\right)^2 + \left((1-a) + \frac{1}{1-a}\right)^2$, for $a \in (0, 1)$. Let $g(x) = \left(x + \frac{1}{x}\right)^2$

---

45

for $x > 0$. Then $f(a) = g(a) + g(1-a)$.

Compute:

$$g(x) = x^2 + 2 + \frac{1}{x^2}, \quad g'(x) = 2x - \frac{2}{x^3}, \quad g''(x) = 2 + \frac{6}{x^4} > 0.$$

Hence, $g$ is strictly convex on $(0, \infty)$.

Since $f''(a) = g''(a) + g''(1-a)$ and each term is positive for $a \in (0, 1)$, we have $f''(a) > 0$. Thus, $f$ is strictly convex on $(0, 1)$ A strictly convex $C^1$-function on an open interval has at most one critical point, which, if it exists, is its global minimum. We compute:

$$f'(a) = g'(a) - g'(1-a) = 0$$

By symmetry, we can solve it and get $a = \frac{1}{2}$. Hence the unique solution of $f'(a) = 0$ in $(0, 1)$ is $a = \frac{1}{2}$.

Therefore, the smallest possible value is $\boxed{C = f\left(\frac{1}{2}\right) = \frac{25}{2}}$.

**Expert Comment:** The initial solution gets the correct value $C = \frac{25}{2}$, but lacks rigor: it incorrectly assumes that symmetry implies a unique minimum, and references the second derivative without computing it (red-highlighted part). The critic accurately points out these gaps. The refined solution corrects them by proving convexity through second derivatives and confirming $a = \frac{1}{2}$ is the unique minimizer through first derivatives (blue-highlighted part). This demonstrates that self-criticism effectively improves the mathematical soundness of the argument.

## C.6 Few-shot Evaluation

We also investigated the effect of few-shot prompting on the INEQMATH test set. Specifically, we compared zero-shot, one-shot, and three-shot configurations across different models.

As shown in Figure 19, the gains in overall accuracy from few-shot prompting were small, typically below 2% compared to zero-shot performance. For instance, Grok 3 achieved 3.5% accuracy in the zero-shot setting but dropped slightly in the one-shot and three-shot settings (2.5% and 1.5%, respectively). Similarly, o1 peaked at 8.0% in both the one-shot and three-shot settings, with minimal difference across shots.
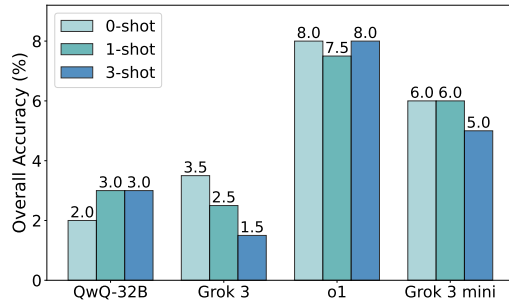


Figure 19: Model performance under zero-shot, one-shot, and three-shot settings (*Overall Accuracy*).
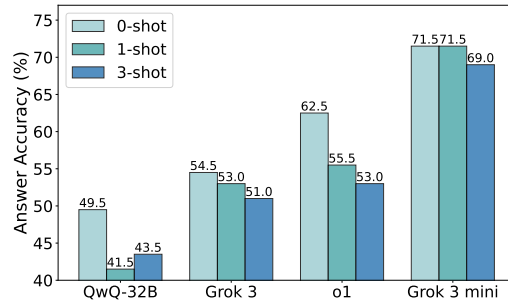
Figure 20: Model performance under zero-shot, one-shot, and three-shot settings (*Answer Accuracy*).

Moreover, Figure 20 shows that few-shot prompting typically reduces answer accuracy. For o1, accuracy drops from 62.5% in the zero-shot setting to 55.5% with one-shot and 53.0% with three-shot. QwQ-32B displays the same trend: both one-shot and three-shot underperform the zero-shot baseline (41.5% and 43.5% vs. 49.5%). These declines suggest overfitting to exemplars, indicating that few-shot prompting does not reliably improve the answer accuracy on the INEQMATH test set.

## C.7 Evaluation on the Formalized INEQMATH

To expand the impact of INEQMATH, we conduct a formal evaluation on state-of-the-art automated theorem proving (ATP) models. The key step in this evaluation is the **formalization process**, which converts the natural language inequality problems in INEQMATH into machine-verifiable Lean4 code.

As illustrated in Figures 21 and Figures 22, this process proceeds in two stages in our experiment. First, we reformulate the inequality problems into proof-style problems using GPT-4.1 [46], ensuring they are structured for formalization. Second, we employ the Goedel-Formalizer-V2-32B [36] to automatically translate these reformulated proof problems into valid Lean4 representations.
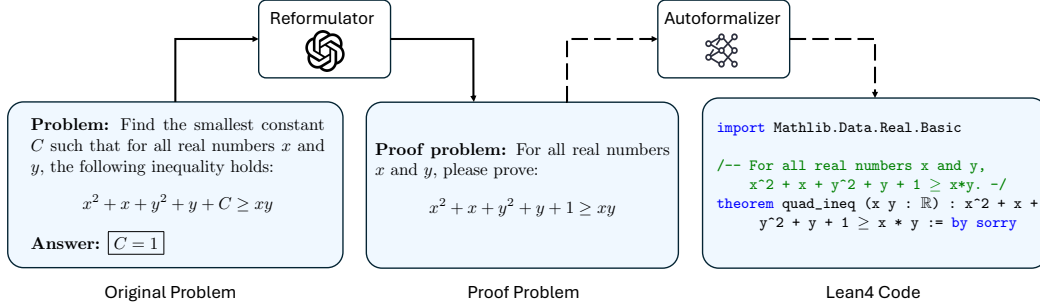


Figure 21: Illustration of the formalization process for bound problems.
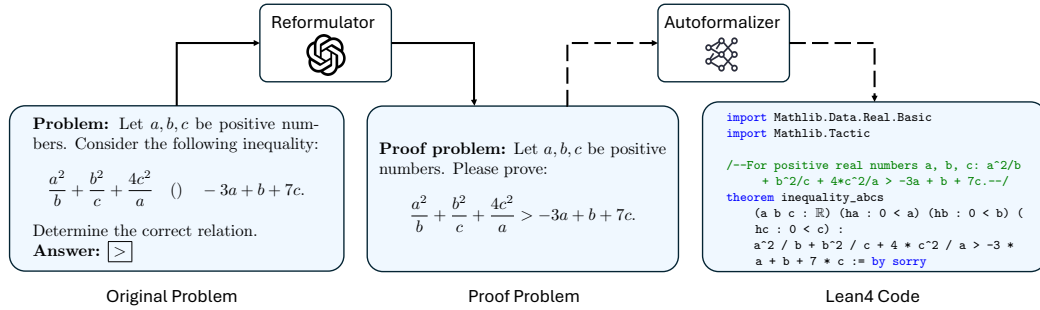


Figure 22: Illustration of the formalization process for relation problems.

Once formalized, we evaluate SOTA ATP models on the Lean4 problems to measure their ability to solve inequality tasks. The results are as follows.

| Model name | Pass rate (Pass@32) |
|---|---|
| DeepSeek-Prover-V2-7B [55] | 6.0% |
| Kimina-Prover-Distill-8B [4] | 12.0% |
| Goedel-Prover-V2-32B [36] | 13.0% |
| Goedel-Prover-SFT [35] | 14.0% |

Table 7: Pass@32 performance of state-of-the-art formal automated theorem proving models.

The results in Table 7 show that state-of-the-art (SOTA) formal automated theorem proving models still suffer from the difficult inequality problems in INEQMATH. Even the best-performing model, Goedel-Prover-SFT, achieves only a 14.0% pass rate, while others remain far lower. This demonstrates that current approaches are inadequate for reliably solving the inequality-focused tasks presented in INEQMATH, and further methods are needed to achieve significant improvements in handling these challenging problems.

## C.8 Memorization Probe

To further demonstrate the modest degree of contamination in INEQMATH, we conducted a memorization probe experiment. In this probe, we systematically rephrased all test problems by swapping the terms on either side of each inequality and then re-evaluated models on the reformulated version. The rephrased problem is mathematically equivalent to the original one, differing only in presentation. This allows us to test whether models had merely memorized the original problems or could generalize to equivalent but rephrased tasks. Examples of the rephrased problems are as follows.

---

**Memorization Probe Reformulation Example 1: Bound Problem**

**Original Problem:** Find the smallest constant $C$ such that for all real numbers $x$ and $y$, the following inequality holds:
$$x^2 + x + y^2 + y + C \geq xy$$
**Original Answer:** $\boxed{C = 1}$

**Rephrased Problem:** Find the smallest constant $C$ such that for all real numbers $x$ and $y$, the following inequality holds:
$$xy \leq x^2 + x + y^2 + y + C$$
**Rephrased Answer:** $\boxed{C = 1}$

---

**Memorization Probe Reformulation Example 2: Relation Problem**

**Original Problem:** Let $a, b, c$ be positive numbers. Consider the following inequality:
$$\frac{a^2}{b} + \frac{b^2}{c} + \frac{4c^2}{a} \quad () \quad -3a + b + 7c.$$
Determine the correct inequality relation to fill in the blank.
Options: (A) $\leq$     (B) $\geq$     (C) $=$     (D) $<$     (E) $>$     (F) None of the above

**Original Answer:** $\boxed{(E) >}$

**Rephrased Problem:** Let $a, b, c$ be positive numbers. Consider the following inequality:
$$-3a + b + 7c \quad () \quad \frac{a^2}{b} + \frac{b^2}{c} + \frac{4c^2}{a}.$$
Determine the correct inequality relation to fill in the blank.
Options: (A) $\leq$     (B) $\geq$     (C) $=$     (D) $<$     (E) $>$     (F) None of the above

**Rephrased Answer:** $\boxed{(D) <}$

---

We evaluate Claude Sonnet 4, GPT-4.1 mini, and o4-mini on both the original and reformulated versions of the INEQMATH test set, with their performance results summarized below.

As shown in Figures 23 and 24, model performance remains largely consistent across the original and reformulated versions of the INEQMATH test set. For example, GPT-4.1 mini maintains an overall accuracy of 8.5% in both conditions, while o4-mini shows only a slight drop from 15.5% to 15.0%. In terms of answer accuracy, Claude Sonnet 4 decreases modestly from 44.0% to 40.0%, whereas o4-mini remains steady at 65.0%. These small shifts—generally under 5 percentage points—indicate that the models adapt well to rephrased tasks rather than relying on memorized solutions. This provides strong evidence that contamination is unlikely, as performance is not driven by rote recall.
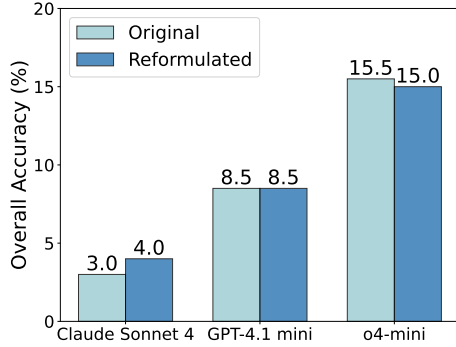
Figure 23: Model performance on the original and reformulated version of the INEQMATH test set (*Overall Accuracy*).
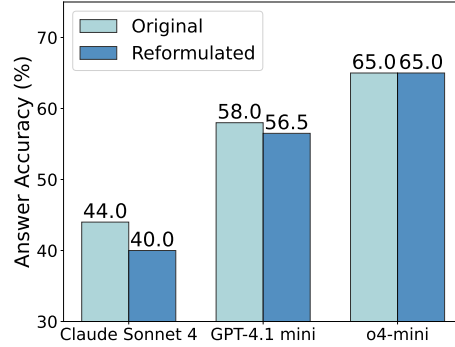
Figure 24: Model performance on the original and reformulated version of the INEQMATH test set (*Answer Accuracy*).

# D Limitations

While our work introduces a novel dataset and evaluation judges for LLM-based inequality proving, we acknowledge several limitations that warrant discussion and offer avenues for future research.

**Potential for data contamination.** Although we took significant measures to mitigate data leakage by commissioning novel test problems curated by experts, keeping ground truth answers private, and utilizing an online leaderboard for evaluation, a residual risk of contamination remains. LLMs possess vast training corpora, and it is possible they have encountered problems with similar structures or underlying principles during pre-training, potentially inflating performance beyond true generalization capabilities. Our expert curation and review process aimed to minimize this, but perfect isolation from prior knowledge is challenging to guarantee.

**Training dataset scale and scope.** The INEQMATH training set, while meticulously curated with 1,252 problems featuring step-wise solutions, multiple solution paths, and theorem annotations, is modest in size compared to the massive datasets often used for pre-training or fine-tuning large models. We prioritized quality and depth (step-wise solutions, theorems) to the challenging Olympiad-level domain over sheer quantity. While sufficient for benchmarking current models, post-training, and exploring test-time techniques, this scale might be insufficient for training highly specialized models from scratch or for capturing the full diversity of inequality types. Future work could focus on scaling up the dataset while maintaining quality, potentially through community contributions.

**Inherent inaccuracies in LLM-as-judge evaluation.** Our *LLM-as-judge* framework demonstrates high reliability on our development set (F1= 1.0 for the *final-answer judge*, > 0.9 average for step-wise judges). However, while significantly more scalable than human expert evaluation, these judges are still imperfect. As illustrated by examples in §B.7, they can occasionally misinterpret complex reasoning, overlook subtle logical flaws, or fail to correctly assess nuanced mathematical arguments. The current set of step-wise judges targets common failure modes but does not cover all possible error types, such as the correctness of complex symbolic transformations or the optimal choice of strategy. Potential improvements include using more powerful (but potentially more expensive) LLMs as judge backends (e.g., o3), developing specialized judges trained on annotated errors, or adding judges for specific mathematical operations like symbolic manipulation verification.

**Mitigation, not elimination, of answer guessability.** The inclusion of step-wise judges significantly mitigates the issue of models guessing the correct final answer without sound reasoning. However, it does not eliminate this possibility entirely. A model might still arrive at the correct bound or relation through chance or heuristics and support it with plausible-sounding, yet flawed, intermediate steps capable of misleading one or more judges. The requirement to pass all judges reduces this risk, but the fundamental challenge of distinguishing genuine mathematical insight from convincing yet spurious reasoning remains.

**Computational cost of evaluation.** While more efficient than manual expert grading, our multi-judge evaluation protocol is computationally more intensive than simple final-answer checking (e.g., string matching). Evaluating each solution requires multiple LLM inferences (one for the final answer, four for step-wise checks). This cost scales linearly with the number of models and problems being evaluated and could become a factor in very large-scale benchmarking efforts.

# E   Broader Impacts

This research focuses on advancing the mathematical reasoning capabilities of LLMs, specifically in the domain of inequality proving. While the work is primarily foundational and unlikely to lead directly to malicious applications such as disinformation or surveillance, potential negative societal impacts could arise from the misuse or misinterpretation of the technology. The most significant risk stems from over-reliance on LLM-generated proofs that may appear correct superficially (achieving high answer accuracy) but contain critical logical flaws, as demonstrated by the sharp drop in performance under our step-wise evaluation. If such flawed proofs were uncritically accepted in fields requiring mathematical rigor, such as scientific modeling, engineering design, or financial analysis, they could lead to incorrect conclusions, faulty systems, or economic miscalculations. Our contribution of a rigorous, step-wise evaluation methodology serves as a potential mitigation strategy by promoting transparency and enabling the identification of fragile reasoning chains, thereby encouraging cautious deployment and emphasizing the need for verification, especially in high-stakes applications. The public release of the INEQMATH benchmark further supports community efforts in understanding and improving the reliability of LLM reasoning.

# NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The abstract and introduction (§1) clearly state the paper's contributions: an informal yet verifiable task formulation for inequality proving (§2), the release of the INEQMATH benchmark with expert curation and annotations (§3), a novel LLM-as-judge evaluation suite (§4), and a systematic empirical study revealing performance gaps and scaling trends (§5). These claims accurately reflect the content presented.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: A dedicated section for limitations (§D) is included in the supplementary material, indicating that limitations are discussed.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.

- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory assumptions and proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

   Answer: [NA]

   Justification: The paper focuses on introducing a new benchmark, evaluation methodology, and empirical analysis of LLM performance on inequality proving. It does not propose new theoretical results that require formal assumptions and proofs.

   Guidelines:

   - The answer NA means that the paper does not include theoretical results.
   - All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
   - All assumptions should be clearly stated or referenced in the statement of any theorems.
   - The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
   - Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
   - Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental result reproducibility**

   Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

   Answer: [Yes]

   Justification: The paper details the dataset (§3), the models evaluated (§5.1), the evaluation methodology including the LLM-as-judge framework and specific judges (§4), and provides prompts and setup details in the appendix (§B, §C). This information should be sufficient for others to understand how the results were obtained and attempt replication. The release of the dataset and leaderboard further aids reproducibility.

   Guidelines:

   - The answer NA means that the paper does not include experiments.

- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

   Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

   Answer: [Yes]

   Justification: Yes, we provide open access to the INEQMATH benchmark via Hugging Face (https://huggingface.co/spaces/OctoTools/octotools), the evaluation code via GitHub (https://github.com/ineqmath/ineqmath), an interactive leaderboard (https://huggingface.co/spaces/OctoTools/octotools), and a project website (https://ineqmath.github.io/). These resources include instructions and are intended to allow faithful reproduction of the main experimental results.

   Guidelines:
   - The answer NA means that paper does not include experiments requiring code.
   - Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
   - While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
   - The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
   - The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
   - The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental setting/details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section §5.1 describes the models, zero-shot setting, and evaluation metrics. The appendix (§C) provides further details, including prompts and placeholders for hyperparameters (e.g., max tokens, temperature), suggesting these details are included. The dataset splits (train/dev/test) are defined in §3.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment statistical significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: No. Results are reported as point estimates (accuracy) on the 200-problem test set. Given the substantial size of the test set and the large performance gaps observed between different model tiers (e.g., >10-20% absolute differences in Table 3), we believe the point estimates provide a clear picture of relative capabilities, though formal statistical significance testing was not performed.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments compute resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Yes. Experiments with proprietary and open-source LLMs were conducted via API calls to their respective providers (e.g., OpenAI, Google, Anthropic, xAI, Together AI), as described in §5.1. For these models, compute infrastructure (e.g., hardware, memory) is managed by the providers and not directly relevant for reproducibility, which depends on the model endpoint and API parameters (§C). Execution time varies with API load and query complexity, while total cost depends on token usage and provider pricing.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code of ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research involves benchmarking LLMs on mathematical tasks and releasing a related dataset. Based on the description, the work appears to align with the NeurIPS Code of Ethics, with no obvious ethical issues concerning data privacy, bias amplification, or malicious use arising from the core methodology or dataset itself.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The paper includes a placeholder for a Broader Impacts section (§E) in the supplementary material, indicating intent to discuss these aspects.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The primary asset released is a dataset of mathematical inequality problems (INEQMATH). While data contamination is a general concern for benchmarks, this dataset does not inherently pose a high risk for direct misuse in the sense intended by the question (e.g., generating harmful content, disinformation). The LLMs evaluated are existing models, and their release is not part of this work.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes. Existing assets used, primarily the evaluated LLMs, are properly credited via citations and mentions (§1, §5.1, Table 4). Their use respects their respective terms of service and copyright principles. For clarity regarding the assets introduced by this paper, the INEQMATH dataset will be released under CC BY-SA 4.0 and the accompanying codebase under the MIT license.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

    Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

    Answer: [Yes]

    Justification: Yes. The paper introduces the new INEQMATH dataset. This asset comprises an expert-curated test benchmark (200 Olympiad-level problems) and a training corpus (1,252 problems). The training data is enriched with step-wise solutions and annotations linking problems to 83 distinct named theorems. Documentation detailing the dataset's curation process, statistical properties (Table 1), format, and examples is provided within the main paper (§3) and the appendix (§A, §A.4), and is intended to be provided alongside the released assets. The release mechanism (e.g., via Hugging Face leaderboard) will contain further usage instructions.

    Guidelines:

    - The answer NA means that the paper does not release new assets.
    - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
    - The paper should discuss whether and how consent was obtained from people whose asset is used.
    - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

    Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

    Answer: [NA]

    Justification: The paper describes using "IMO-level medalists" and "experts" for dataset curation (§3). This appears to be expert consultation or commissioned work rather than crowdsourcing or research involving human subjects in the typical sense requiring detailed reporting of instructions and compensation within the paper.

    Guidelines:

    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
    - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
    - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

    Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

    Answer: [NA]

    Justification: The research methodology, involving expert mathematicians for problem creation/review and evaluating LLMs, does not constitute human subjects research that would typically require IRB approval or involve risks to participants.

    Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: LMs are central to this research. They are the primary subjects of evaluation (§5). Furthermore, LLMs are used in a non-standard way as part of the core methodology itself: for data rephrasing during curation (§3) and as automated judges in the evaluation framework (§4). This usage is explicitly described.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.