DEEPTHEOREM: ADVANCING LLM REASONING FOR THEOREM PROVING THROUGH NATURAL LANGUAGE AND REINFORCEMENT LEARNING

Anonymous authorsPaper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

033 034

035

037

040

041

042

043

044

046 047

048

051

052

ABSTRACT

Theorem proving serves as a major testbed for evaluating complex reasoning abilities in large language models (LLMs). However, traditional automated theorem proving (ATP) approaches rely heavily on formal proof systems that poorly align with LLMs' strength derived from informal, natural language knowledge acquired during pre-training. To fully leverage the theorem-proving knowledge acquired from pre-training, in this work, we present DeepTheorem, a comprehensive informal theorem-proving suite exploiting natural language to enhance LLM mathematical reasoning. DeepTheorem includes 1) a large-scale dataset of 121K high-quality IMO-level informal theorems and proofs spanning diverse mathematical domains, rigorously annotated for correctness, difficulty, and topic categories, accompanied by systematically constructed verifiable theorem variants; 2) adaptation of RL-Zero explicitly to informal theorem proving, leveraging the verified theorem variants to incentivize robust mathematical inference; 3) comprehensive outcome and process evaluation metrics examining proof correctness and the quality of reasoning steps; and 4) a novel informal theorem proving benchmark consoliadted from three established math competitions, formatted for automatic evaluation. Extensive experimental analyses demonstrate DeepTheorem significantly improves LLM theorem-proving performance compared to existing datasets and supervised fine-tuning protocols, achieving state-of-the-art accuracy and reasoning quality. Our findings highlight DeepTheorem's potential to fundamentally advance automated informal theorem proving and mathematical exploration.

1 Introduction

Theorem proving is widely regarded as a pinnacle challenge for evaluating advanced reasoning capabilities of both human and artificial intelligence. It requires integrating diverse cognitive facets such as abstraction, strategic inference, pattern recognition, and meticulous logical deduction. Recent advancements in deep learning, especially in large language models (LLMs), have significantly reshaped the landscape of automated theorem proving (ATP). Much prior work attempts ATP by integrating LLMs with either formal proof engines such as Lean, Coq, and Isabelle (Zheng et al., 2022; Liu et al., 2023; Tsoukalas et al., 2024) or domain-specific languages from ProofWiki (Welleck et al., 2022). However, these proof methods impose a significant barrier for LLMs whose primary strength derives from the vast corpus of natural language and LaTeX-based mathematical texts used during pre-training. This inherent misalignment limits LLMs' capability in theorem proving, leaving a considerable gap between their potential and actual performance.

In this paper, we present **DeepTheorem**, a novel, comprehensive suite expressly designed to leverage natural language to unleash the latent mathematical reasoning ability of LLMs for theorem proving. Instead of relying on formal proof assistants, DeepTheorem offers a scalable, intuitive, and flexible alternative, enabling LLMs to generate informal mathematical proofs aligning closely with human mathematicians' heuristic-driven thinking. Central to our approach is the construction of a large-scale benchmark comprising 121K IMO-level informal mathematical theorems with precise annotations such as correctness labels, difficulty levels, diverse mathematical domains, and verifiable theorem variants amenable to advanced reinforcement learning paradigms. Additionally,

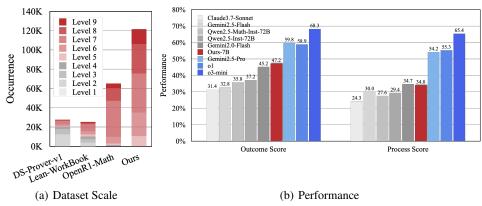


Figure 1: (a): Our dataset surpasses others with extremely challenging theories; (b): RL-Zero training with our DeepTheorem datasets on 7B model achieves strong results.

we adapt reinforcement learning algorithms to the problem of informal theorem proving for the first time, significantly enhancing LLM's ability to reason mathematically beyond the constraints of supervised fine-tuning (SFT). To evaluate existing and newly trained models on informal theorem proving, we also construct a new benchmark sourced from established mathematics competitions, and propose comprehensive evaluation metrics that rigorously assess the correctness of generated proofs and the processes underlying the proofs themselves.

Through extensive experiments, we show that leading LLMs still exhibit significant limitations in theorem proving. However, when trained with DeepTheorem, they achieve substantial performance improvements over models trained using existing datasets, showcasing the effectiveness of our natural-language-focused approach. Our results underscore the promise of DeepTheorem to redefine LLM-driven mathematical reasoning, offering a robust platform for continued progress in automated and scalable informal theorem proving.

To sum up, our key contributions are:

- We introduce the *DeepTheorem* framework, a comprehensive informal theorem-proving suite exploiting natural language to enhance LLM mathematical reasoning;
- We open-source a large-scale natural-language theorem collection of 121K informal mathematical theorems and corresponding high-quality proofs at IMO-level difficulty, suitable for both SFT and RL;
- We innovatively adapt the RL-Zero training method explicitly to informal theorem proving, significantly enhancing LLM's reasoning capacity beyond traditional SFT methods;
- We introduce a new benchmark for evaluating informal theorem proving, and develop a comprehensive evaluation framework assessing both the correctness of theorem proofs (outcome evaluation) and the completeness, logical validity, and correctness of generated reasoning processes (process evaluation);
- Through extensive experiments, we establish the superiority of our DeepTheorem training paradigm, achieving state-of-the-art performance and surpassing existing informal theorem datasets and training methods;

2 Dataset

Overview The *DeepTheorem* dataset¹ is a novel, large-scale resource designed to advance LLMs in informal mathematical theorem reasoning. Mined from a diverse web corpus, it addresses the need for challenging, decontaminated, and diverse topics to push LLMs toward frontier theorem proving. As illustrated in Figure 2, each entry in the dataset offers distinct features tailored to support diverse research objectives, including: 1) a mathematical theorem in standardized LaTeX format; 2) a True-or-False correctness label, where a False label indicates that the theorem can be mathematically

¹In the rest of this section, *DeepTheorem* dataset refers specifically to the DeepTheorem training dataset.

```
Question: Let f(z) = \frac{e^{iz}}{z}, and consider the contour integral \int_{\mathbb{R}} f(z) dz over the real line. Assume the contour is closed by a semicircular arc in the upper half-plane. Prove or disprove that the value of the improper integral \int_0^\infty \frac{\sin x}{x} dx = \Im\left(\int_0^\infty \frac{e^{ix}}{x} dx\right)

Because the integrand \frac{e^{iz}}{z} is meromorphic with a simple pole at z = 0, we consider the principal value of the integral along the whole real axis. That is, we consider

Final Answer: True

Difficulty: 8

Topic:
Calculus -> Integral Calculus -> Techniques Of Integration Calculus -> Integral Calculus -> Applications Of Integrals Algebra -> Intermediate Algebra -> Complex Numbers
```

Figure 2: A data sample from the DeepTheorem dataset.

disproved; 3) a fine-grained difficulty score ranging from level $5\sim10$; 4) mathematical topics for targeted analysis; and 5) detailed step-by-step proof solutions generated by o3-mini.

Comprising approximately 121K curated samples, *DeepTheorem* empowers large-scale exploration of LLM mathematical reasoning, curriculum learning, and cross-domain generalization. It serves as a versatile resource for advancing automated theorem proving, enhancing model reasoning capability, and developing adaptive learning frameworks, positioning it as a cornerstone for next-generation NLP research in mathematical theorem proving and logical reasoning. The key features of *DeepTheorem* are:

- Large scale: Unlike smaller training corpus such as Lean-Workbook (Ying et al., 2024), Deepseek-Prover-v1's training corpus (Xin et al., 2024), and theorems from OpenR1-Math (Face, 2025), *DeepTheorem* dataset leverages the vastness of web-sourced content to ensure comprehensive coverage of mathematical concepts and problem types. As shown in Figure 1(a), our dataset consists of approximately 121K theorems, significantly outscaling prior datasets.
- Frontier and extremely challenging theorems: DeepTheorem dataset is distinguished by its inclusion of advanced mathematical theorems, each annotated with difficulty levels to enable targeted evaluation and training across a spectrum of complexities. As shown in Figure 4, DeepTheorem dataset emphasizes theorems at high difficulty levels (6–9), surpassing existing corpora in complexity and challenge, presenting significant



Figure 3: Statistics of *DeepTheorem* dataset hierarchical topics.

challenges for state-of-the-art LLMs while aligning with frontier, IMO-level benchmarks such as FIMO (Liu et al., 2023).

- **Diverse topics**: As shown in Figure 3, *DeepTheorem* dataset captures the breadth of informal theorem-based reasoning by covering nearly the entirety of the mathematical landscape, including algebra, discrete math, applied math, calculus, geometry, mathematical analysis, number theory, etc. By encompassing this wide array of domains, *DeepTheorem* dataset enables researchers to assess model performance on both specialized and interdisciplinary mathematical tasks, fostering the development of LLMs that can generalize effectively across the full spectrum of mathematical reasoning.
- Strict decontamination: To preserve evaluation integrity, DeepTheorem dataset employs rigorous decontamination processes to avoid overlap with widely used benchmarks. The targets of our decontamination includes general math reasoning benchmarks MATH (Hendrycks et al.,

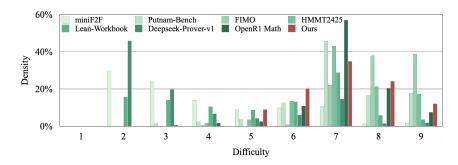


Figure 4: Difficulty density histogram comparison across widely used training dataset (i.e. Lean-Workbook, Deepseek-Prover-V1, OpenR1 Math) and testing benchmarks (i.e. miniF2F, Putnam-Bench, FIMO, HMMT2425).

2021b), AIME (MAA, a), AMC (MAA, b), Minerva Math (Lewkowycz et al., 2022), Olympiad-Bench (He et al., 2024), Omni-MATH (Gao et al., 2025), MathOdyssey (Fang et al., 2024), GAOKAO (Zhong et al., 2024), JEEBench (Arora et al., 2023), MMLU-STEM (Hendrycks et al., 2021a), CMATH (Wei et al., 2023), OlympicArena (Huang et al., 2024), GSM8K (Cobbe et al., 2021), GPQA (Rein et al., 2024) - and theorem proving benchmarks: miniF2F (Zheng et al., 2022), PutnamBench (Tsoukalas et al., 2024), FIMO (Liu et al., 2023), and HMMT (Harvard-MIT Mathematics Tournament, 2024, 2025).

• **Proofs from advanced LLMs**: *DeepTheorem* dataset includes concise, high-quality proof solutions generated by o3-mini, tailored for supervised fine-tuning (SFT). These proofs provide a compact yet complete outline of the logical steps required to prove (or disprove) each theorem, optimized for clarity and brevity. Unlike verbose or overly formal proofs, these proofs, expressed in LaTeX, align with the informal nature of LLMs, making them an effective learning signal. By incorporating these proofs, the dataset enables models to internalize structured reasoning patterns, improving their ability to generate coherent and logically sound mathematical arguments.

2.1 Constructing DeepTheorem Dataset

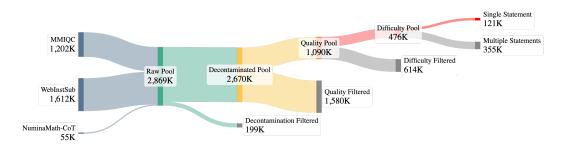


Figure 5: The process pipeline of constructing *DeepTheorem* dataset.

As shown in Figure 5, the construction pipeline of *DeepTheorem* dataset starts by aggregating raw data from multiple sources, including MMIQC (Liu et al., 2025), WebInstruct (Yue et al., 2024), and NuminaMath-CoT (LI et al., 2024).

Decontamination To ensure the integrity of *DeepTheorem* dataset, we implement a rigorous decontamination process to eliminate overlap between training and testing datasets. The process employs a recall-and-justify pipeline to identify and remove potential duplicates, consisting of three key steps:

1. *Embedding Generation*: We use an embedding model² (Reimers & Gurevych, 2019; Toshniwal et al., 2025) to generate sentence embeddings for all theorem statements in the training and testing datasets.

²paraphrase-multilingual-MiniLM-L12-v2

- 2. *Similarity Recall*: For each training sample, we compute its embedding cosine similarity to all test samples, and recall the top five testing samples exceeding a similarity threshold (set to 0.7).
- 3. *Contamination Justification*: An LLM³ evaluates whether the recalled test samples are contaminated within the current training sample (Section D.1).

This process removed approximately 199K contaminated samples, effectively identifying identical cases, generalized questions, and converse theorems, Removed contaminated examples are shown in Appendix A. Approximately 2.6M samples remain for the next processing step.

Quality control and proof generation We also implement a rigorous quality control pipeline for generating and validating theorem statements and their proofs. The process involves four key steps:

- 1. *Theorem Justification*: An LLM verifies that the question is complete, and indeed a theorem-proving question (Section D.2).
- 2. *Rationale Summarization*: An LLM summarizes the original question and generates a formatted, concise, self-contained theorem (Section D.3).
- 3. *Proof Generation*: o3-mini (high effort) generates the proof solution with True-or-False conclusion about the theorem (Section D.4).
- 4. *Logical Validation*: The LLM performs an extra justification step to check that the theorem-proof pair is logically coherent. (Section D.5)

This systematic approach yielded 1.08M high-quality, mathematically sound theorem-proof pairs.

Difficulty and single statement annotation We annotate the difficulty levels of *DeepTheorem* dataset, and remove questions with multiple statements to prove.

- 1. Difficulty Annotation: An LLM analyzes each theorem statement following the strategy of Gao et al. (2025), considering factors such as logical complexity, mathematical prerequisites, and proof length, to assign a difficulty score on a scale of 1 to 9. Only questions with a difficulty score of at least 5 are retained. (Section D.6)
- Single-Statement Filtering: We filter out samples that query for proving multiple statements, retaining only those with a single, well-defined theorem to ensure clarity and consistency with evaluation.
- 3. *Topic Annotation*: Finally, we annotate the topic domain of the mathematical theorems with LLMs. (Section D.7)

The difficulty and single statements filtering results in 121K challenging theorems, yielding the final *DeepTheorem* dataset.

3 THEOREM PROVING VIA REINFORCEMENT LEARNING

Motivation Conventionally, informal theorem-proving datasets are utilized through supervised fine-tuning (SFT), where models learn to generate proofs by imitating dataset examples. However, recent studies on RL-Zero demonstrate its superior performance over SFT by leveraging a base model's pretrained knowledge and exploratory capabilities (Jaech et al., 2024; DeepSeek-AI et al., 2025). This raises a natural question: *Can we harness the base model's exploration ability for informal theorem proving?* In this section, we explore the possibility of utilizing RL-Zero for informal theorem proving. The process involves three key steps: 1) data augmentation to generate contradictory theorem variants for binary rewards; 2) RL-Zero training with GRPO (Shao et al., 2024); and 3) Evaluation of the theorem-proof generation.

3.1 Theorems with Verifiable Rewards

Theorems can be disproved To construct a theorem with rewards for RL-Zero, we make the key observation:

³ GPT-40 is used for annotation in this section unless otherwise specified.

A statement need not be correct but can be also proven incorrect, enabling a binary reward structure compatible with RL-Zero.

This observation allows us to transform *DeepTheorem*'s theorems into true-or-false variants, facilitating RL training that incentivizes robust reasoning.

To construct such training data, we use an LLM to expand the original theorems into contradictory variants that can be *disproved*. Specifically, we strictly limit the transformation made to the original theorem, so that the resulting variant is either entailed by or contradictory from the original theorem. Consider the example in Table 1 (omitting the hypotheses for simplicity): if the original theorem can be proved, Variant #1 is also correct and can be mathematically proved in the same manner as the original one, while Variant #2 must be incorrect and can be disproved.

Table 1: An example of theorem variants given an original theorem.

Theorems	Example
Original	x > 1
Variant 1	x > 0
Variant 2	x < 1

With such logically entailing or contradictory transformations, we are able to construct variants of a theorem that are guaranteed to be correct or incorrect by only accessing the theorem itself but not the proof process, which makes this transformation task much easier than annotating new math statements, and thus allowing a relatively weaker LLM (e.g. Qwen2.5-72B-Instruct, Yang et al., 2024) to perform it. After this expansion phase, we further annotate the completeness of the resulting theorem pool and finally acquire a training set of 242K mathematical theorems that can either be proved or disproved, each with a complete proof trajectory (see Appendix D.8 for more details).

3.2 BINARY REWARDS ACTIVATE THEOREM PROOF GENERATION

With the aforementioned theorem variants, we can now apply reinforcement learning to natural language theorem proving. Specifically, we adopt the GRPO algorithm (Shao et al., 2024).

Proof generation with RL Inspired by the success of reasoning-specialized models such as R1 and its open-source reproductions (DeepSeek-AI et al., 2025; Hu et al., 2025), we encourage the model to enclose its reasoning process in <think> </think> tags in the system prompt to incentivize more detailed reasoning behaviours (see Appendix D.9), and then ask the model to end each proof with either "\boxed{proved}" or "\boxed{disproved}". In the reward function, we extract this answer and compare it against the ground truth, giving a reward of 1 if the answer matches, and 0 otherwise. We also enforce several sanity checks to prevent model collapse: if the ratio of white spaces in a model's solution is less than 0.05 or the average character repetition count is greater than 300, then a reward of 0 is issued regardless of the answer.

3.3 EVALUATION

The theorem-proving questions used for evaluation are drawn from two challenging benchmarks — FIMO (Liu et al., 2023) and Putnam (Tsoukalas et al., 2024) — and a newly constructed theorem-proving subset of HMMT (Harvard-MIT Mathematics Tournament, 2024, 2025).

Outcome evaluation Evaluating the correctness of natural language (NL) proofs poses a significant challenge, as it mirrors the complexity faced by humans in assessing the logical coherence and mathematical validity of informal reasoning. Unlike formal theorem-proving systems that rely on structured logic, NL proofs lack a standardized format, making their evaluation inherently subjective and difficult to automate. To address this, we propose a novel evaluation framework that leverages multiple en-

Table 2: Test data statistics. Each original theorem is manually expanded into multiple entailing or contradictory variants. *Random accuracy* indicates the expected score of random guessing following the outcome criteria described below.

Bench	Scale	Variants (Avg.)	Random Acc.
FIMO	172	2.7	17.4
HMMT	205	3.5	11.2
Putnam	281	2.9	15.4

tailing and contradictory variants derived from each theorem. By assessing the model's ability to consistently assign correct truth values across these variants, we indirectly estimate its theorem justification ability. When the number of variants is sufficiently large, this approach provides a robust proxy for evaluating the correctness of NL proof generation.

Thus, we manually expand each question in the three data sources into multiple entailing or contradictory variants following the same variant generation protocol in Section 3.1, and the resulting benchmarks are shown in Table 2. When evaluating a model, we ask it to either prove or disprove each theorem and corresponding variants, and evaluate the results with the criteria below:

Outcome Criteria

A test case in a theorem testing set is passed if and only if:

- 1. The model explicitly produces a truth value (true or false) for theorems and variants;
- 2. The predicted truth value for the original theorem is correct;
- 3. The predicted truth values for all entailing variants are the same as the original theorem;
- 4. The predicted truth values for all contradictory variants are the inverse of the original theorem.

Process evaluation Since theorem proving requires generating logically validated proofs for each reasoning step, we also develop a process evaluation framework that evaluates the quality of proof along four dimensions:

- Logical Validity: Check if each step follows logically from the previous one. Flag any logical errors.:
- Completeness: Verify if all necessary cases and steps are included to prove the theorem;
- Correctness: Confirm if the final conclusion is correct;
- Clarity: Assess if the proof is clear, unambiguous, and well-explained.

We use GPT-40 as the LLM judge and ask it to score the proof using a weighted sum of the four dimensions (prompt given in Appendix D.10). In Appendix B, we also present the results using o3-mini as the judge as well as human evaluation.

4 EXPERIMENTS

4.1 SETTINGS

We train two sets of models, using supervised fine-tuning (SFT) and zero reinforcement learning (RL-Zero) respectively, both starting from Qwen2.5-Base (Yang et al., 2024). For SFT, we train the models for 3 epochs on the complete proof solutions in the dataset, using one machine for training each model. For RL-Zero, we adopt GRPO with batch size 128, group size 64, and maximum rollout length 8192. We train the models for 1000 steps, and distribute each model across two machines during training. Following the settings of Hu et al. (2025), we do not apply any KL regularization or entropy loss, as we find that KL regularization has a negligible impact on model performance, while entropy loss leads to model collapse.

As a baseline, we select the theorem-proving subset of OpenR1-Math (Face, 2025), the highest-quality existing theorem-proving dataset with complete questions and responses. We apply the same processing pipelines to it as detailed in Section 2.1, which yields 66K original theorems and 130K variants in total. We dub this processed dataset *OpenR1-Math-Proof*.

4.2 MAIN RESULTS

DeepTheorem with RL-Zero achieves the best performance The main results are presented in Table 3. *DeepTheorem* demonstrates superior performance over OpenR1-Math-Proof, especially

Table 3: Outcome (out.) and Process (proc.) evaluation of models trained on OpenR1-Math-Proof and DeepTheorem.

Model	Strategy	Data	FIN	МО	IO HM		Putnam		Avg.	
	~		out.	proc.	out.	proc.	out.	proc.	out.	proc.
1.5B	SFT	OpenR1-Proof DeepTheorem	20.63 31.75	8.66 18.86	11.86 15.25	4.80 9.41	35.42 36.46	18.98 21.43	22.64 27.82	10.81 16.57
RL	OpenR1-Proof DeepTheorem	34.92 31.75	8.54 15.23	16.95 23.73	5.10 10.15	55.21 52.08	17.92 22.79	35.69 35.85	10.52 16.06	
2D	SFT	OpenR1-Proof DeepTheorem	23.81 33.33	12.85 20.38	15.25 20.34	6.90 12.15	43.75 36.46	27.96 25.43	27.60 30.04	15.90 19.32
3B	RL	OpenR1-Proof DeepTheorem	34.92 38.10	14.33 23.39	23.73 25.42	11.72 13.56	57.29 52.08	35.11 33.84	38.65 38.53	20.39 23.60
	SFT	OpenR1-Proof DeepTheorem	30.16 34.92	18.23 26.69	15.25 22.03	8.63 15.41	48.96 41.67	32.95 33.50	31.46 32.87	19.94 25.20
7B	RL	OpenR1-Proof DeepTheorem	42.86 55.56	22.79 39.07	25.42 28.81	13.15 20.85	60.42 57.29	38.94 42.20	42.90 47.22	24.96 34.04

for the 7B backbone and in terms of process evaluation. On the other hand, our RL-Zero training paradigm consistently outperforms SFT, validating the effectiveness of RL-Zero in pushing the models' reasoning capabilities beyond the limit of SFT.

DeepTheorem achieves strong parameter efficiency We demonstrates that our DeepTheorem-RL strategy achieves strong parameter efficiency in Figure 6. Compared to the Qwen2.5 series, training DeepTheorem on 1.5 to 7B models significantly improves the informal theorem proving boundary at parameter-performance space. Moreover, when extrapolated DeepTheorem parameter efficiency also surpasses SOTA commercial models such as o1 and o3-mini.

SOTA performance at equal model scale In Table 4, we also provide the evaluation results of SOTA LLMs on the three benchmarks. These results suggest that theorem proving, especially our newly constructed HMMT benchmark, is still quite challenging for LLMs. On

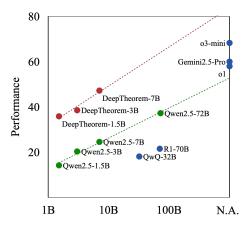


Figure 6: Performance improvement of models trained with DeepTheorem over baselines on theorem proving benchmarks.

the other hand, our 7B model, trained with RL-Zero on *DeepTheorem*, outperforms SOTA models of much larger sizes, including those specialized in math and reasoning, demonstrating the superior quality of *DeepTheorem* and our innovative outcome-supervised RL training approach for theorem proving.

Reasoning with theorem proving skills In Figure 7, we visualize the techniques used by our 7B model trained with RL on *DeepTheorem*, where direct proof is most commonly used, followed by proof by exhaustion and construction. In Appendix C, we provide a non-cherry-picked example generation, finding the model to deliver a clear and correct disproof, highlighting its efficacy in tackling advanced mathematical problems with precision and clarity.

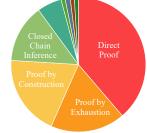


Figure 7: Distribution of proof techniques used by DeepTheorem-7B.

5 RELATED WORK

Theorem proving Theorem proving is a non-trivial task for natural language processing. In the context of LLMs, theorem proving is typically achieved via formal languages such as Lean, Coq,

Table 4: Evaluation comparison of our model trained with RL on DeepTheorem with SOTA LLMs for both commercial models and open source models. *: denotes our method. Inst: Instruct; DS: Deepseek.

Model	FIN	МО	HMMT		Putnam		Avg.		#Rank	
	out.	proc.	out.	proc.	out.	proc.	out.	proc.	out.	proc.
Commercial Models										
Claude3.7-Sonnet	34.92	26.28	13.56	8.29	45.83	38.33	31.44	24.30	9	9
Gemini2.5-Pro	57.14	54.06	57.63	49.82	64.58	58.75	59.78	54.21	2	3
Gemini2.5-Flash	30.16	28.95	25.42	22.02	42.71	38.98	32.76	29.98	8	6
GPT-4o	34.92	30.70	16.95	14.59	22.92	18.88	24.93	21.39	12	10
o1-mini	60.32	55.23	35.59	30.90	61.46	52.88	52.46	46.34	4	4
o1	66.67	61.00	47.46	47.30	62.50	57.55	58.88	55.28	3	2
o3-mini	80.95	77.61	45.76	43.47	78.12	75.12	68.28	65.40	1	1
			Open Sc	ource Mo	dels					
Qwen2.5-Inst-7B	30.16	21.13	10.17	6.83	33.33	25.39	24.55	17.78	13	12
Qwen2.5-Inst-72B	49.21	37.35	13.56	9.78	48.96	41.00	37.24	29.38	6	7
Qwen2.5-Math-Inst-7B	28.57	18.86	3.39	1.61	25.00	18.79	18.99	13.09	16	17
Qwen2.5-Math-Inst-72B	47.62	36.02	11.86	8.61	47.92	38.04	35.80	27.56	7	8
DS-Prover-v1.5-RL-7B	25.40	13.81	11.86	6.32	34.38	22.42	23.88	14.18	14	16
DS-Prover-v2-7B	30.16	21.86	5.08	1.71	40.62	28.54	25.29	17.37	11	13
R1-Distill-7B	6.35	4.27	0	0	4.17	2.58	3.51	2.28	18	18
R1-Distill-70B	17.46	14.05	16.95	13.52	30.21	23.10	21.54	16.89	15	14
QwQ-32B	17.46	15.41	11.86	10.10	25.00	18.19	18.11	14.57	17	15
Llama3.3-Inst-70B	41.27	27.33	10.17	4.12	36.46	25.30	29.30	18.92	10	11
*DeepTheorem-RL-7B	55.56	39.07	28.81	20.85	57.29	42.20	47.22	34.04	5	5

and Isabelle (Zheng et al., 2022; Liu et al., 2023; Tsoukalas et al., 2024). LLMs specialized at theorem proving have been proposed over the years, represented by GPT-f (Polu & Sutskever, 2020), DeepSeek-Prover (Xin et al., 2024; 2025; Ren et al., 2025), TheoremLlama (Wang et al., 2024), InterLM-StepProver (Wu et al., 2024), MPS-Prover (Liang et al., 2025), Goedel-Prover (Lin et al., 2025), and Kimina-Prover (Wang et al., 2025). Although most of these works focus on formal theorem proving, NaturalProofs (Welleck et al., 2021) and NaturalProver (Welleck et al., 2022) have emerged as some of the few works that attend to informal theorem proving. However, NaturalProofs use a domain-specific language from the ProofWiki website. In comparison, *DeepTheorem* represents the first attempt at exploring the more commonly used LaTeX-based natural language theorem proving at scale.

Learning to reason with RL-zero RL-Zero (DeepSeek-AI et al., 2025) is a streamlined framework designed to develop reinforcement learning capabilities in LLMs without SFT. While recent advances in LLM reasoning have been significantly influenced by RL techniques (Jaech et al., 2024; DeepSeek-AI et al., 2025; Team, 2024; xAI, 2025; Google, 2025), existing approaches predominantly focus on closed-form questions, addressing only a subset of reasoning problems. In contrast, we investigate the application of RL-Zero in process-oriented reasoning, specifically in informal theorem-proving. To the best of our knowledge, this is the first study to apply RL-Zero to informal theorem proving, marking a significant advance in enabling LLMs to address more diversified reasoning tasks in mathematical and logical domains.

6 Conclusion

In this paper, we introduce DeepTheorem, a novel comprehensive theorem-proving suite involving a large-scale annotated dataset of 121K IMO-level informal mathematical theorems and corresponding high-quality natural-language proofs, alongside systematically constructed verifiable theorem variants. We further adapt RL-Zero method to informal theorem reasoning, significantly surpassing supervised fine-tuning in performance. Comprehensive evaluations involving outcome accuracy and detailed process assessment on our newly constructed benchmark demonstrate the effectiveness of our approach, achieving state-of-the-art theorem-proving performance and significantly pushing LLM reasoning boundaries. Through these contributions, DeepTheorem provides a robust foundation for future advancements in automated mathematical theorem proving, leveraging natural language flexibility to empower scalable, human-like reasoning abilities in large language models.

REFERENCES

- Daman Arora, Himanshu Singh, and Mausam. Have LLMs advanced enough? a challenging problem solving benchmark for large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 7527–7543, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.468. URL https://aclanthology.org/2023.emnlp-main.468.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168, 2021. URL https://arxiv.org/abs/2110.14168.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. CoRR, abs/2501.12948, 2025. doi: 10. 48550/ARXIV.2501.12948. URL https://doi.org/10.48550/arXiv.2501.12948.
- Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL https://github.com/huggingface/open-r1.
- Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and Kai Zou. Mathodyssey: Benchmarking mathematical problem-solving skills in large language models using odyssey math data. *CoRR*, abs/2406.18321, 2024. doi: 10.48550/ARXIV.2406.18321. URL https://doi.org/10.48550/arXiv.2406.18321.
- 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, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=yaqPf0KAlN.
- Google. Gemini 2.0 flash thinking, 2025. URL https://cloud.google.com/vertex-ai/generative-ai/docs/thinking. Accessed on March 25, 2025.
- Harvard-MIT Mathematics Tournament. Hmmt dataset and resources. https://www.hmmt.org/, 2024, 2025. Accessed: 2025-05-12.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. Olympiadbench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 3828–3850. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.211. URL https://doi.org/10.18653/v1/2024.acl-long.211.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021a. URL https://openreview.net/forum?id=d7KBjmI3GmQ.

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 Joaquin Vanschoren and Sai-Kit Yeung (eds.), Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual, 2021b. URL https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html.

Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model. *CoRR*, abs/2503.24290, 2025. doi: 10.48550/ARXIV.2503.24290. URL https://doi.org/10.48550/arXiv.2503.24290.

Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, Yiwei Qin, Yan Ma, Jiadi Su, Yixiu Liu, Yuxiang Zheng, Shaoting Zhang, Dahua Lin, Yu Qiao, and Pengfei Liu. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent AI. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/222d2eaf24cf8259a35d6c7130d31425-Abstract-Datasets_and_Benchmarks_Track.html.

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, and Ilge Akkaya. Openai o1 system card. CoRR, abs/2412.16720, 2024. doi: 10.48550/ARXIV.2412.16720. URL https://doi.org/10.48550/arXiv.2412.16720.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 3843–3857. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/18abbeef8cfe9203fdf9053c9c4fe191-Paper-Conference.pdf.

Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. Numinamath. [https://huggingface.

```
co/AI-MO/NuminaMath-1.5] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.
```

- Zhenwen Liang, Linfeng Song, Yang Li, Tao Yang, Feng Zhang, Haitao Mi, and Dong Yu. Mpsprover: Advancing stepwise theorem proving by multi-perspective search and data curation, 2025. URL https://arxiv.org/abs/2505.10962.
- Yong Lin, Shange Tang, Bohan Lyu, Jiayun Wu, Hongzhou Lin, Kaiyu Yang, Jia Li, Mengzhou Xia, Danqi Chen, Sanjeev Arora, and Chi Jin. Goedel-prover: A frontier model for open-source automated theorem proving. *CoRR*, abs/2502.07640, 2025. doi: 10.48550/ARXIV.2502.07640. URL https://doi.org/10.48550/arXiv.2502.07640.
- Chengwu Liu, Jianhao Shen, Huajian Xin, Zhengying Liu, Ye Yuan, Haiming Wang, Wei Ju, Chuanyang Zheng, Yichun Yin, Lin Li, Ming Zhang, and Qun Liu. FIMO: A challenge formal dataset for automated theorem proving. *CoRR*, abs/2309.04295, 2023. doi: 10.48550/ARXIV. 2309.04295. URL https://doi.org/10.48550/arxiv.2309.04295.
- Haoxiong Liu, Yifan Zhang, Yifan Luo, and Andrew C. Yao. Augmenting math word problems via iterative question composing. In Toby Walsh, Julie Shah, and Zico Kolter (eds.), *AAAI-25*, *Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 March 4, 2025, Philadelphia, PA, USA*, pp. 24605–24613. AAAI Press, 2025. doi: 10.1609/AAAI. V39I23.34640. URL https://doi.org/10.1609/aaai.v39i23.34640.
- MAA. American invitational mathematics examination (AIME). Mathematics Competition Series, a. URL https://maa.org/math-competitions/aime.
- MAA. American mathematics competitions (AMC 10/12). Mathematics Competition Series, b. URL https://maa.org/math-competitions/amc.
- Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *CoRR*, abs/2009.03393, 2020. URL https://arxiv.org/abs/2009.03393.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bertnetworks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pp. 3980–3990. Association for Computational Linguistics, 2019. doi: 10.18653/V1/D19-1410. URL https://doi.org/10.18653/v1/D19-1410.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. URL https://openreview.net/forum?id=Ti67584b98.
- 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. *CoRR*, abs/2504.21801, 2025. URL https://arxiv.org/abs/2504.21801.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024. doi: 10.48550/ARXIV.2402.03300. URL https://doi.org/10.48550/arXiv.2402.03300.
- Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown, Nov 2024. URL https://qwenlm.github.io/blog/qwq-32b-preview/.
- Shubham Toshniwal, Wei Du, Ivan Moshkov, Branislav Kisacanin, Alexan Ayrapetyan, and Igor Gitman. Openmathinstruct-2: Accelerating AI for math with massive open-source instruction data. In *The Thirteenth International Conference on Learning Representations, ICLR* 2025, Singapore, April 24-28, 2025. OpenReview.net, 2025. URL https://openreview.net/forum?id=mTCbq2QssD.

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. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/1582eaf9e0cf349e1e5a6ee453100aa1-Abstract-Datasets_and_Benchmarks_Track.html.

- 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 Saxcé, 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. *CoRR*, abs/2504.11354, 2025. doi: 10.48550/ARXIV.2504.11354. URL https://doi.org/10.48550/arxiv.2504.11354.
- Ruida Wang, Jipeng Zhang, Yizhen Jia, Rui Pan, Shizhe Diao, Renjie Pi, and Tong Zhang. Theoremllama: Transforming general-purpose llms into lean4 experts. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pp. 11953–11974. Association for Computational Linguistics, 2024. URL https://aclanthology.org/2024.emnlp-main.667.
- Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. CMATH: can your language model pass chinese elementary school math test? *CoRR*, abs/2306.16636, 2023. doi: 10.48550/ARXIV. 2306.16636. URL https://doi.org/10.48550/arXiv.2306.16636.
- Sean Welleck, Jiacheng Liu, Ronan Le Bras, Hanna Hajishirzi, Yejin Choi, and Kyunghyun Cho. Naturalproofs: Mathematical theorem proving in natural language. In Joaquin Vanschoren and Sai-Kit Yeung (eds.), Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual, 2021. URL https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/d9d4f495e875a2e075a1a4a6e1b9770f-Abstract-round1.html.
- Sean Welleck, Jiacheng Liu, Ximing Lu, Hannaneh Hajishirzi, and Yejin Choi. Natural-prover: Grounded mathematical proof generation with language models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/1fc548a8243ad06616eee731e0572927-Abstract-Conference.html.
- 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. *CoRR*, abs/2410.15700, 2024. doi: 10.48550/ARXIV.2410.15700. URL https://doi.org/10.48550/arXiv.2410.15700.
- xAI. Grok: Artificial intelligence assistant, 2025. URL https://x.ai. Developed by xAI, accessed on March 25, 2025.
- Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data. *CoRR*, abs/2405.14333, 2024. doi: 10.48550/ARXIV.2405.14333. URL https://doi.org/10.48550/arXiv.2405.14333.
- Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong Shao, Wanjia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Haowei Zhang, Qihao Zhu, Dejian Yang, Zhibin Gou,

Z. F. Wu, Fuli Luo, and Chong Ruan. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=I4YAIwrsXa.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *CoRR*, abs/2412.15115, 2024. doi: 10.48550/ARXIV.2412.15115. URL https://doi.org/10.48550/arXiv.2412.15115.

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 Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/bf236666a2cc5f3ae05d2e08485efc4c-Abstract-Datasets_and_Benchmarks_Track.html.

Xiang Yue, Tianyu Zheng, Ge Zhang, and Wenhu Chen. Mammoth2: Scaling instructions from the web. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/a4ca07aa108036f80cbb5b82285fd4b1-Abstract-Conference.html.

Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. minif2f: a cross-system benchmark for formal olympiad-level mathematics. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. URL https://openreview.net/forum?id=9ZPegFuFTFv.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024, pp. 2299–2314. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-NAACL.149. URL https://doi.org/10.18653/v1/2024.findings-naacl.149.

A EXAMPLES OF DECONTAMINATED TRAINING CASES

Table 5: Examples of benchmark contamination in polynomial and number theory problems. Generalizing and logically equivalent parts are highlighted.

Contaminated Example	Benchmark Example	Relationship
Let $p(x)$ be a univariate polynomial. Then $p(x)$ is nonnegative for all $x \in \mathbb{R}$ if and only if $p(x)$ can be expressed as a sum of squares (SOS), i.e., $p(x) = \sum_{i=1}^k q_i^2(x)$ for some polynomials $q_1(x), \ldots, q_k(x)$.	Let $p(x)$ be a polynomial that is nonnegative for all real x . Prove that for some k , there are polynomials $f_1(x), \ldots, f_k(x)$ such that $p(x) = \sum_{j=1}^k (f_j(x))^2$.	Identical
Let $p(x_1, x_2,, x_n)$ be a real polynomial. If $p(x_1, x_2,, x_n)$ is non-negative for all $(x_1, x_2,, x_n) \in \mathbb{R}^n$, then $p(x_1, x_2,, x_n)$ can be expressed as a sum of squares of polynomials if and only if p belongs to the quadratic module generated by the constraints of a certain semialgebraic set. Formally, there exists a set of polynomials $q_i(x_1, x_2,, x_n)$ such that $p(x_1, x_2,, x_n) = \sum_{i=1}^k q_i(x_1, x_2,, x_n)^2$, provided certain conditions on p and the domain hold to ensure the SOS representation.	Let $p(x)$ be a polynomial that is nonnegative for all real x . Prove that for some k , there are polynomials $f_1(x), \ldots, f_k(x)$ such that $p(x) = \sum_{j=1}^k (f_j(x))^2$.	Generalizing
Let n be a positive integer. If n is not prime, then $2^n - 1$ is not prime.	Show that if n is a positive integer and $2^n - 1$ is prime, then n is prime.	Logically Converse

B COMPARISON OF PROCESS EVALUATION JUDGES

B.1 COMPARISON BETWEEN GPT-40 AND 03-MIN

In Table 6, we present the comparison between using GPT-40 and o3-mini as judges for process evaluation. While o3-mini is stricter and gives lower scores on average, the relative ranks of all evaluated models are similar, and the scores between the two judges have a correlation coefficient of more than 0.95, demonstrating strong consistency.

Table 6: Comparison of process evaluation scores using different judges.

Model	GPT	Г-4о	o3-mini		
Model	Score	Rank	Score	Rank	
Claude3.7-Sonnet	24.30	9	15.54	8	
Gemini2.5-Pro	54.21	3	53.98	2	
Gemini2.5-Flash	29.98	6	29.42	5	
GPT-4o	21.39	10	14.29	12	
o1-mini	46.34	4	36.02	4	
o1	55.28	2	53.52	3	
o3-mini	65.40	1	65.57	1	
Qwen2.5-Inst-7B	17.78	12	8.03	15	
Qwen2.5-Inst-72B	29.38	7	15.31	9	
Qwen2.5-Math-Inst-7B	13.09	17	5.17	16	
Qwen2.5-Math-Inst-72B	27.56	8	15.08	10	
DS-Prover-v1.5-RL-7B	14.18	16	4.34	17	
DS-Prover-v2-7B	17.37	13	8.36	14	
R1-Distill-7B	2.28	18	1.90	18	
R1-Distill-70B	16.89	14	16.46	6	
QwQ-32B	14.57	15	14.89	11	
Llama3.3-Inst-70B	18.92	11	10.05	13	
DeepTheorem-RL-7B	34.04	5	15.95	7	

B.2 HUMAN EVALUATION

To further validate our LLM-based process evaluation frameworks, we also conduct human evaluations on a subset of the model outputs. Due to the mathematical difficulty of the questions involved, we sampled 10 questions and their corresponding proofs generated by three representative models: o3-mini, R1-Distill-70B, and Llama3.3-Instruct-70B, covering both API-based and open-source models, as well as reasoning-specialized and general-purpose models.

We asked two human experts (both PhD in computer science) to conduct pairwise evaluations of the model responses, and the win rates are presented in Table 7. The scores from o3-mini in Table 6 are also converted into pairwise binary scores by comparing the scores assigned to responses from two models. The results suggest that the evaluation of the LLM judge aligns well with human experts, resulting in consistent rankings.

Table 7: Comparison of pairwise evaluations from human experts and o3-mini.

	Н	Human Evaluation o3-mini Evaluation					ion	
Model	Win			Rank			Lose	Rank
o3-mini	0.75	0.00	0.25	1	0.67	0.17	0.17	1
Llama3.3-Inst-70B	0.17	0.08	0.75	3	0.00	0.17	0.83	3
R1-Distill-70B	0.50	0.06	0.44	2	0.50	0.25	0.25	2

C EXAMPLE OF MODEL-GENERATED PROOF

Model-generated proof

Question:

Let P be a polynomial of degree n>1 with integer coefficients and let k be any positive integer. Consider the polynomial $Q(x)=P(P(\dots P(P(x))\dots))$, with k pairs of parentheses. Prove or disprove that Q has more than n integer fixed points, i.e. integers satisfying the equation Q(x)=x.

Model Solution:

<think>

To determine whether the polynomial $Q(x) = P(P(\dots P(P(x)) \dots))$ with k pairs of parentheses has more than n integer fixed points, we need to analyze the properties of the polynomial P and how the composition of P affects its fixed points.

First, recall that P(x) is a polynomial of degree n>1 with integer coefficients. The number of integer roots of P(x)-x=0 (i.e., the number of integer fixed points of P) is at most n by the fundamental theorem of algebra and properties of polynomials. Let's denote the integer fixed points of P by r_1, r_2, \ldots, r_m , where $m \leq n$.

Now, consider the polynomial $Q(x) = P(P(\dots P(P(x)) \dots))$ with k pairs of parentheses. We need to determine the number of integer solutions to the equation Q(x) = x. If r is an integer fixed point of Q, then r must be an integer fixed point of some composition of P. Specifically, if r is a fixed point of Q, there exists some $m \leq k$ such that $P^{(m)}(r) = r$, where $P^{(m)}$ denotes the m-th iterate of P.

The key observation is that the number of integer fixed points of Q is bounded by the number of integer preperiodic points of P. Since P is a polynomial of degree n, the number of integer preperiodic points of P is at most n for each iterate. However, we need to consider the composition Q which is the k-th iterate of P.

The polynomial Q(x) is also of degree n^k , and the number of integer roots of Q(x) - x = 0 is at most n^k . However, we need to check if it is possible for Q(x) to have more than n integer fixed points. Since each application of P can introduce at most n new fixed points, the total number of integer fixed points of Q is still bounded by the number of preperiodic points of P, which is at most n for each iterate.

To be more precise, let's consider the fixed points. If r is a fixed point of Q, then r must be a root of the polynomial $P^k(x)-x$. The polynomial $P^k(x)-x$ is of degree n^k , and by the fundamental theorem of algebra, it can have at most n^k roots. However, we are interested in integer roots, and the number of integer roots is at most the number of integer preperiodic points of P, which is at most n for each iterate.

Therefore, the number of integer fixed points of Q(x) is at most n, the same as the number of integer fixed points of P for each iterate. Hence, Q(x) cannot have more than n integer fixed points.

Thus, the statement that Q has more than n integer fixed points is false.

<

This is a non-cherry-picked example from our *DeepTheorem*-RL-7B model. We analyze this proof in the following dimensions:

• Correctness. The model accurately identifies that integer fixed points of Q(x) satisfy Q(x) = x, and recognizes that Q(x) has degree n^k , with Q(x) - x = 0 having at most n^k roots. Its insight that the number of integer fixed points is constrained by the fixed points of P, which are at most

n, effectively guides the solution. The model leverages the iterative structure of Q(x) to argue that additional compositions do not increase the number of integer fixed points beyond n, aligning with expected polynomial behavior.

- Clarity. The proof is well-structured, progressing logically from problem definition to conclusion. Mathematical notation, such as $P^{(m)}$ for the m-th iterate, is precise and standard, enhancing readability. The solution is concise, making it accessible and suitable for a mathematical audience
- **Strengths.** The model's ability to distill a complex problem into a clear argument showcases its strength in handling polynomial compositions. Its focus on fixed points and iterative properties demonstrates a solid grasp of the underlying mathematics, making it a valuable tool for such proofs.

In summary, the model delivers a clear and correct disproof, highlighting its efficacy in tackling advanced mathematical problems with precision and clarity.

D DeepTheorem PROMPT TEMPLATE

D.1 CONTAMINATION JUSTIFICATION

972

973 974

975976

977 978

979

980

981

982

983

984

985 986 987

988 989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1008 1009

1010

1011 1012

1013

1014

1015

1016

1017

1018

1021

1023

1024 1025

Meta Prompt for Contamination Justification

I will now give you two theorems: the Original Theorem and the Candidate Theorem. Please help me determine if the following two theorems are the same.

Original Theorem: ***theorem1***

Candidate Theorem: ***theorem2***

Disregard the names and minor changes in word order. If their theorem prompts are very similar, without considering the proving process, we consider them to be the same theorem. Note that you should not consider the solution process, only the theorem prompts.

You should only respond with True or False. Do not respond with anything else.

D.2 THEOREM-PROVING ANNOTATION

Meta Prompt for Theorem-Proving Annotation

You are an expert in classifying questions based on their type and intent. Given the following discussion:

Discussion

- problem: ***problem***
- solution: ***solution***

Determine whether the question is:

- A question-answering (QA) question seeking a specific value or factual response, or
- A theorem-proving question requiring logical reasoning, derivation, or proof of a mathematical or theoretical statement.

Provide a clear classification (QA or theorem-proving) and justify your decision with a concise explanation. Consider the following:

- QA questions typically ask for a specific fact, value, or definitive answer (e.g., "What is the capital of France?" or "What is the value of x in 2x = 8?").
- Theorem-proving questions typically involve logical reasoning, mathematical derivation, or proving a general statement (e.g., "Prove that the sum of two even numbers is even" or "Derive the Pythagorean theorem").

Return "True" if the question is a theorem-proving question, and "False" if it is a QA question.

D.3 RATIONALE SUMMARIZATION

Meta Prompt for Rationale Summarization

You are provided with a corpus of forum discussions about mathematical topics. Your task is to analyze the discussion and:

- 1. Identify the key mathematical concepts, ideas, or rationales driving the discussion.
- 2. Act as a teacher to formulate a theorem based on the discussion, presented as a formal theorem statement.
- # Requirements
- All mathematical equations must be formatted in LaTeX.
- The theorem should be a clear, formal statement (e.g., "Let $f: \mathbb{R} \to \mathbb{C}$ be a smooth function, ...").
- The output must be in JSON format, with the following structure:

"rationale": "A description of the main mathematical concepts or ideas in the discussion.", "theorem": "A formal theorem statement based on the discussion."

Discussion

```
- problem: ***problem***
- solution: ***solution***
```

D.4 PROOF GENERATION

1026

1027

103310341035

1036 1037

1039

1040

1041

1042 1043

1044 1045

1046 1047

1048

1049

1050

1057 1058

1059

1061

1062

1063

1064

1065

1067

1068

1069

1070

1071

1074

1075

1076

1077

1078

1079

Meta Prompt for Proof Generation from o3-mini (high)

You are provided with a corpus of forum discussions about mathematical topics.

A theorem statement is summarized from the discussion. Your task is to provide a proof for the theorem statement based on the discussion.

Requirements

- All mathematical equations must be formatted in LaTeX.
- The proof should be a clear, formal statement (e.g., "To prove this theorem, we can start by ...").
- The output must be in JSON format, with the following structure:

```
"proof": "A proof for the theorem statement."
```

- # Discussion
- problem: ***problem***
- solution: ***solution***
- # Theorem Rationale
- ***theorem***

D.5 LOGICAL VALIDATION

Meta Prompt for Logical Validation of the Proof

You are an expert in mathematical theorem proving and logical analysis. Given the following theorem and its proof or disproof, your task is to analyze each step of the proof or disproof to determine if it is valid, providing a detailed justification for each step's correctness or identifying any errors.

```
# Theorem
```

- ***theorem***
- # Proof or Disproof
- ***Proof***
- # Instructions
- 1. **Analyze Each Step**:
- Verify if the step is mathematically correct, logically sound, and relevant to proving or disproving the theorem.
- Check for adherence to mathematical definitions, theorems, or properties cited in the step.
- Ensure the step follows from previous steps or given assumptions without logical gaps.
- If the step involves a disproof, confirm that it correctly demonstrates a counterexample or contradiction.
- 2. **Overall Assessment**:
- Conclude whether the entire proof or disproof is valid.
- If invalid, return False and summarize the critical errors and recommend how to fix the proof/disproof.
- If valid, return True and confirm that it fully addresses the theorem.

D.6 DIFFICULTY ANNOTATION

Meta Prompt for Difficulty Annotation

CONTEXT

I am a teacher, and I have some high-level olympiad math problems.

I want to evaluate the difficulty of these math problems. There are some references available regarding the difficulty of the problems:

<difficulty reference>

Examples for difficulty levels For reference, here are problems from each of the difficulty levels 1-10:

- 1: How many integer values of x satisfy $|x| < 3\pi$? (2021 Spring AMC 10B, Problem 1)
- 1.5: A number is called flippy if its digits alternate between two distinct digits. For example, 2020 and 37373 are flippy, but 3883 and 123123 are not. How many five-digit flippy numbers are divisible by 15? (2020 AMC 8, Problem 19)
- 2: A fair 6-sided die is repeatedly rolled until an odd number appears. What is the probability that every even number appears at least once before the first occurrence of an odd number? (2021 Spring AMC 10B, Problem 18)
- 2.5: A, B, C are three piles of rocks. The mean weight of the rocks in A is 40 pounds, the mean weight of the rocks in B is 50 pounds, the mean weight of the rocks in the combined piles A and B is 43 pounds, and the mean weight of the rocks in the combined piles A and C is 44 pounds. What is the greatest possible integer value for the mean in pounds of the rocks in the combined piles B and C? (2013 AMC 12A, Problem 16)
- 3: Triangle ABC with $\overline{AB} = 50$ and $\overline{AC} = 10$ has area 120. Let D be the midpoint of \overline{AB} , and let E be the midpoint of \overline{AC} . The angle bisector of $\angle BAC$ intersects \overline{DE} and \overline{BC} at F and G, respectively. What is the area of quadrilateral FDBG? (2018 AMC 10A, Problem 24)
- 3.5: Find the number of integer values of k in the closed interval [-500, 500] for which the equation $\log(kx) = 2\log(x+2)$ has exactly one real solution. (2017 AIME II, Problem 7) 4: Define a sequence recursively by $x_0 = 5$ and

$$x_{n+1} = \frac{x_n^2 + 5x_n + 4}{x_n + 6}$$

for all nonnegative integers n. Let m be the least positive integer such that

$$x_m \le 4 + \frac{1}{2^{20}}.$$

In which of the following intervals does m lie?

- (A) [9,26] (B) [27,80] (C) [81,242] (D) [243,728] (E) $[729,\infty)$ (2019 AMC 10B, Problem 24 and 2019 AMC 12B, Problem 22)
- 4.5: Find, with proof, all positive integers n for which $2^n + 12^n + 2011^n$ is a perfect square. (USAJMO 2011/1)
- 5: Find all triples (a, b, c) of real numbers such that the following system holds:

$$a+b+c = \frac{1}{a} + \frac{1}{b} + \frac{1}{c},$$

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

(JBMO 2020/1)

- 5.5: Triangle ABC has $\angle BAC = 60^{\circ}$, $\angle CBA \le 90^{\circ}$, BC = 1, and $AC \ge AB$. Let H, I, and O be the orthocenter, incenter, and circumcenter of $\triangle ABC$, respectively. Assume that the area of pentagon BCOIH is the maximum possible. What is $\angle CBA$? (2011 AMC 12A, Problem 25)
- 6: Let $\triangle ABC$ be an acute triangle with circumcircle ω , and let H be the intersection of the altitudes of $\triangle ABC$. Suppose the tangent to the circumcircle of $\triangle HBC$ at H intersects ω at points X and Y with HA=3, HX=2, and HY=6. The area of $\triangle ABC$ can be written

in the form $m\sqrt{n}$, where m and n are positive integers, and n is not divisible by the square of any prime. Find m+n. (2020 AIME I, Problem 15)

6.5: Rectangles BCC_1B_2 , CAA_1C_2 , and ABB_1A_2 are erected outside an acute triangle ABC. Suppose that

$$\angle BC_1C + \angle CA_1A + \angle AB_1B = 180^{\circ}.$$

Prove that lines B_1C_2 , C_1A_2 , and A_1B_2 are concurrent. (USAMO 2021/1, USAJMO 2021/2)

7: We say that a finite set S in the plane is balanced if, for any two different points A, B in S, there is a point C in S such that AC = BC. We say that S is centre-free if for any three points A, B, C in S, there is no point P in S such that PA = PB = PC. Show that for all integers $n \geq 3$, there exists a balanced set consisting of n points. Determine all integers $n \geq 3$ for which there exists a balanced centre-free set consisting of n points. (IMO 2015/1) 7.5: Let \mathbb{Z} be the set of integers. Find all functions $f: \mathbb{Z} \to \mathbb{Z}$ such that

$$xf(2f(y) - x) + y^{2}f(2x - f(y)) = \frac{f(x)^{2}}{x} + f(yf(y))$$

for all $x, y \in \mathbb{Z}$ with $x \neq 0$. (USAMO 2014/2)

8: For each positive integer n, the Bank of Cape Town issues coins of denomination $\frac{1}{n}$. Given a finite collection of such coins (of not necessarily different denominations) with total value at most most $99 + \frac{1}{2}$, prove that it is possible to split this collection into 100 or fewer groups, such that each group has total value at most 1. (IMO 2014/5)

8.5: Let I be the incentre of acute triangle ABC with $AB \neq AC$. The incircle ω of ABC is tangent to sides BC, CA, and AB at D, E, and F, respectively. The line through D perpendicular to EF meets ω at R. Line AR meets ω again at P. The circumcircles of triangle PCE and PBF meet again at Q. Prove that lines DI and PQ meet on the line through A perpendicular to AI. (IMO 2019/6)

9: Let k be a positive integer and let S be a finite set of odd prime numbers. Prove that there is at most one way (up to rotation and reflection) to place the elements of S around the circle such that the product of any two neighbors is of the form $x^2 + x + k$ for some positive integer x. (IMO 2022/3)

9.5: An anti-Pascal triangle is an equilateral triangular array of numbers such that, except for the numbers in the bottom row, each number is the absolute value of the difference of the two numbers immediately below it. For example, the following is an anti-Pascal triangle with four rows which contains every integer from 1 to 10.

$$\begin{array}{rrr}
 & 4 \\
 & 2 & 6 \\
 & 5 & 7 & 1 \\
 & 8 & 3 & 10 & 9
\end{array}$$

Does there exist an anti-Pascal triangle with 2018 rows which contains every integer from 1 to $1+2+3+\cdots+2018$? (IMO 2018/3)

10: Prove that there exists a positive constant c such that the following statement is true: Consider an integer n>1, and a set $\mathcal S$ of n points in the plane such that the distance between any two different points in $\mathcal S$ is at least 1. It follows that there is a line ℓ separating $\mathcal S$ such that the distance from any point of $\mathcal S$ to ℓ is at least $cn^{-1/3}$.

Some known difficulty ratings of the competitions. ### HMMT (November) Individual Round, Problem 6-8: 4

Individual Round, Problem 6-8: 4 Individual Round, Problem 10: 4.5 Team Round: 4-5

Guts: 3.5-5.25 ### CEMC **Part A: 1-1.5** How many different 3-digit whole numbers can be formed using the digits 4, 7, and 9, assuming that no digit can be repeated in a number? (2015 Gauss 7 Problem 10)

Part B: 1-2

Two lines with slopes $\frac{1}{4}$ and $\frac{5}{4}$ intersect at (1,1). What is the area of the triangle formed by these two lines and the vertical line x = 5? (2017 Cayley Problem 19) Part C (Gauss/Pascal):

2-2.5

Suppose that $\frac{2009}{2014} + \frac{2019}{n} = \frac{a}{h}$, where a, b, and n are positive integers with $\frac{a}{h}$ in lowest terms.

Suppose that $\frac{2009}{2014} + \frac{2019}{n} = \frac{a}{b}$, where a, b, and n are positive integers with $\frac{a}{b}$ in lowest terms. What is the sum of the digits of the smallest positive integer n for which a is a multiple of 1004? (2014 Pascal Problem 25)

Part C (Cayley/Fermat): 2.5-3

Wayne has 3 green buckets, 3 red buckets, 3 blue buckets, and 3 yellow buckets. He randomly distributes 4 hockey pucks among the green buckets, with each puck equally likely to be put in each bucket. Similarly, he distributes 3 pucks among the red buckets, 2 pucks among the blue buckets, and 1 puck among the yellow buckets. Once he is finished, what is the probability that a green bucket contains more pucks than each of the other 11 buckets? (2018 Fermat Problem 24)

Indonesia MO

Problem 1/5: 3.5 In a drawer, there are at most 2009 balls, some of them are white, the rest are blue, which are randomly distributed. If two balls were taken at the same time, then the probability that the balls are both blue or both white is $\frac{1}{2}$. Determine the maximum amount of white balls in the drawer, such that the probability statement is true?

Problem 2/6: 4.5 Find the lowest possible values from the function

$$f(x) = x^{2008} - 2x^{2007} + 3x^{2006} - 4x^{2005} + 5x^{2004} - \dots - 2006x^3 + 2007x^2 - 2008x + 2009$$

for any real numbers x.

Problem 3/7: 5 A pair of integers (m, n) is called good if

$$m \mid n^2 + n \text{ and } n \mid m^2 + m$$

Given 2 positive integers a, b > 1 which are relatively prime, prove that there exists a good pair (m, n) with $a \mid m$ and $b \mid n$, but $a \nmid n$ and $b \nmid m$.

Problem 4/8: 6 Given an acute triangle ABC. The incircle of triangle ABC touches BC, CA, AB respectively at D, E, F. The angle bisector of $\angle A$ cuts DE and DF respectively at K and L. Suppose AA_1 is one of the altitudes of triangle ABC, and M be the midpoint of BC. (a) Prove that BK and CL are perpendicular with the angle bisector of $\angle BAC$. (b) Show that A_1KML is a cyclic quadrilateral.

JBMO

Problem 1: 4 Find all real numbers a, b, c, d such that

$$a + b + c + d = 20$$
, $ab + ac + ad + bc + bd + cd = 150$.

Problem 2: 4.5-5 Let ABCD be a convex quadrilateral with $\angle DAC = \angle BDC = 36^\circ$, $\angle CBD = 18^\circ$ and $\angle BAC = 72^\circ$. The diagonals intersect at point P. Determine the measure of $\angle APD$.

Problem 3: 5 Find all prime numbers p, q, r, such that $\frac{p}{q} - \frac{4}{r+1} = 1$.

Problem 4: 6A 4×4 table is divided into 16 white unit square cells. Two cells are called neighbors if they share a common side. A move consists in choosing a cell and changing the colors of neighbors from white to black or from black to white. After exactly n moves all the 16 cells were black. Find all possible values of n.

Problem 1/4: 5 There are a+b bowls arranged in a row, numbered 1 through a+b, where a and b are given positive integers. Initially, each of the first a bowls contains an apple, and each of the last b bowls contains a pear. A legal move consists of moving an apple from bowl i to bowl i+1 and a pear from bowl j to bowl j-1, provided that the difference i-j is even. We permit multiple fruits in the same bowl at the same time. The goal is to end up with the first b bowls each containing a pear and the last a bowls each containing an apple. Show that this is possible if and only if the product ab is even.

Problem 2/5: 6-6.5 Let a,b,c be positive real numbers such that $a+b+c=4\sqrt[3]{abc}$. Prove that

$$2(ab + bc + ca) + 4\min(a^2, b^2, c^2) \ge a^2 + b^2 + c^2.$$

Problem 3/6: 7 Two rational numbers $\frac{m}{n}$ and $\frac{n}{m}$ are written on a blackboard, where m and n are relatively prime positive integers. At any point, Evan may pick two of the numbers x and y written on the board and write either their arithmetic mean $\frac{x+y}{2}$ or their harmonic mean $\frac{2xy}{x+y}$ on the board as well. Find all pairs (m,n) such that Evan can write 1 on the board in finitely many steps.

HMMT (February) Individual Round, Problem 1-5: 5 Individual Round, Problem 6-10: 5.5-6 Team Round: 7.5 HMIC: 8

APMO Problem 1: 6 Problem 2: 7 Problem 3: 7 Problem 4: 7.5 Problem 5: 8.5 ### Balkan MO Problem 1: 5 Solve the equation $3^x - 5^y = z^2$ in positive integers. Problem 2: 6.5 Let MN be a line parallel to the side BC of a triangle ABC, with M on the side AB and N on the side AC. The lines BN and CM meet at point P. The circumcircles of triangles BMP and CNP meet at two distinct points P and Q. Prove that $\angle BAQ = \angle CAP$. Problem 3: 7.5 A 9 × 12 rectangle is partitioned into unit squares. The centers of all the unit squares, except for the four corner squares and eight squares sharing a common side with one of them, are coloured red. Is it possible to label these red centres $C_1, C_2, ..., C_{96}$ in such way that the following to conditions are both fulfilled (i) the distances $C_1C_2, ...C_{95}C_{96}, C_{96}C_1$ are all equal to $\sqrt{13}$ (ii) the closed broken line $C_1C_2...C_{96}C_1$ has a centre of symmetry? Problem 4: 8 Denote by S the set of all positive integers. Find all functions $f: S \to S$ such that

$$f(f^{2}(m) + 2f^{2}(n)) = m^{2} + 2n^{2} \text{ for all } m, n \in S.$$

USAMO Problem 1/4: 6-7 Problem 2/5: 7-8 Three nonnegative real numbers r_1 , r_2 , r_3 are written on a blackboard. These numbers have the property that there exist integers a_1 , a_2 , a_3 , not all zero, satisfying $a_1r_1 + a_2r_2 + a_3r_3 = 0$. We are permitted to perform the following operation: find two numbers x, y on the blackboard with $x \le y$, then erase y and write y - x in its place. Prove that after a finite number of such operations, we can end up with at least one 0 on the blackboard. Problem 3/6: 8-9 Prove that any monic polynomial (a polynomial with leading coefficient 1) of degree n with real coefficients is the average of two monic polynomials of degree n with n real roots.

USA TST Problem 1/4/7: 6.5-7 Problem 2/5/8: 7.5-8 Problem 3/6/9: 8.5-9 ### Putnam Problem A/B,1-2: 7 Find the least possible area of a concave set in the 7-D plane that intersects both branches of the hyperparabola xyz=1 and both branches of the hyperbola xwy=-1. (A set S in the plane is called convex if for any two points in S the line segment connecting them is contained in S.) Problem A/B,3-4: 8 Let S be an S matrix all of whose entries are S and whose rows are mutually orthogonal. Suppose S has an S submatrix whose entries are all 1. Show that S has an S has an S has an another such that S has an another such that S has an another such that S has another such that S has an another such that S has a such that

$$2m = a^{19} + b^{99} + k * 2^{1000}.$$

Problem 2/5: 9 Given a positive integer n=1 and real numbers $a_1 < a_2 < \ldots < a_n$, such that $\frac{1}{a_1} + \frac{1}{a_2} + \ldots + \frac{1}{a_n} \leq 1$, prove that for any positive real number x,

$$\left(\frac{1}{a_1^2+x}+\frac{1}{a_2^2+x}+\ldots+\frac{1}{a_n^2+x}\right)^2 \ge \frac{1}{2a_1(a_1-1)+2x}.$$

Problem 3/6: 9.5-10 Let n>1 be an integer and let a_0,a_1,\ldots,a_n be non-negative real numbers. Define $S_k=\sum_{i=0}^k \binom{k}{i}a_i$ for $k=0,1,\ldots,n$. Prove that

$$\frac{1}{n} \sum_{k=0}^{n-1} S_k^2 - \frac{1}{n^2} \left(\sum_{k=0}^n S_k \right)^2 \le \frac{4}{45} (S_n - S_0)^2.$$

IMO **Problem 1/4: 5.5-7** Let Γ be the circumcircle of acute triangle ABC. Points D and E are on segments AB and AC respectively such that AD = AE. The perpendicular bisectors of BD and CE intersect minor arcs AB and AC of Γ at points F and G respectively. Prove that lines DE and FG are either parallel or they are the same line.

```
1296
             **Problem 2/5: 7-8** Let P(x) be a polynomial of degree n > 1 with integer coefficients,
1297
             and let k be a positive integer. Consider the polynomial Q(x) = P(P(\dots P(P(x)) \dots)),
1298
             where P occurs k times. Prove that there are at most n integers t such that Q(t) = t.
1299
             **Problem 3/6: 9-10** Let ABC be an equilateral triangle. Let A_1, B_1, C_1 be interior
             points of ABC such that BA_1 = A_1C, CB_1 = B_1A, AC_1 = C_1B, and
1300
1301
                                       \angle BA_1C + \angle CB_1A + \angle AC_1B = 480^{\circ}
1302
             Let BC_1 and CB_1 meet at A_2, let CA_1 and AC_1 meet at B_2, and let AB_1 and BA_1 meet
1303
             at C_2. Prove that if triangle A_1B_1C_1 is scalene, then the three circumcircles of triangles
1304
             AA_1A_2, BB_1B_2 and CC_1C_2 all pass through two common points.
1305
             ### IMO Shortlist
             Problem 1-2: 5.5-7
             Problem 3-4: 7-8
             Problem 5+: 9-10
             </difficulty reference>
1309
1310
1311
            # OBJECTIVE #
1312
             1. Summarize the math problem in a brief sentence, describing the concepts involved in the
1313
            math problem.
1314
            2. Based on the source of the given problem, as well as the difficulty of the problems
1315
            referenced in these materials and the solution to the current problem, please provide an
1316
            overall difficulty score for the current problem. The score should be a number between 1
1317
            and 10, with increments of 0.5, and should align perfectly with the materials.
1318
             # STYLE #
             Data report.
1319
             # TONE #
1320
            Professional, scientific.
1321
            # AUDIENCE #
1322
            Students. Enable them to better understand the difficulty of the math problems.
1323
            # RESPONSE: MARKDOWN REPORT #
1324
             ## Summarization
             [Summarize the math problem in a brief paragraph.]
1326
             ## Difficulty
1327
             [Rate the difficulty of the math problem and give the reason.]
1328
            # ATTENTION #
             - Add "=== report over ===" at the end of the report.
1330
1331
             <example math problem>
1332
             [Question]:
1333
             If \frac{1}{9} + \frac{1}{18} = \frac{1}{x}, what is the number that replaces the x to make the equation true?
1334
             [Solution]:
1335
             We simplify the left side and express it as a fraction with numerator 1: \frac{1}{9} + \frac{1}{18} = \frac{2}{18} + \frac{1}{18} = \frac{2}{18}
1336
                =\frac{1}{6}. Therefore, the number that replaces the \square is 6.
1337
             [Source]: 2010_Pascal
1338
             </example math problem>
1339
             ## Summarization
1340
             The problem requires finding a value that makes the equation \frac{1}{9} + \frac{1}{18} = \frac{1}{\square}. This involves
1341
            adding two fractions and determining the equivalent fraction.
            ## Difficulty
            Rating: 1
             Reason: This problem is straightforward and primarily involves basic fraction addition, mak-
            ing it suitable for early middle school students.
1345
             === report over ===
1347
1348
             <example math problem>
1349
             [Question]:
```

Let $\mathcal P$ be a convex polygon with n sides, $n\geq 3$. Any set of n-3 diagonals of $\mathcal P$ that do not intersect in the interior of the polygon determine a triangulation of $\mathcal P$ into n-2 triangles. If $\mathcal P$ is regular and there is a triangulation of $\mathcal P$ consisting of only isosceles triangles, find all the possible values of n.

[Solution]:

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1363

1365

1367

1369

13701371

1372

1373

1374

1375

1376

1380

1382

1384

1386

1387

1388

1389

1390

1391

1392

1393

1394

1399

1400

1401

1402

1403

We label the vertices of \mathcal{P} as $P_0, P_1, P_2, \ldots, P_n$. Consider a diagonal $d = \overline{P_a P_{a+k}}, k \leq n/2$ in the triangulation. We show that k must have the form 2^m for some nonnegative integer m. This diagonal partitions \mathcal{P} into two regions \mathcal{Q}, \mathcal{R} , and is the side of an isosceles triangle in both regions. Without loss of generality suppose the area of Q is less than the area of R (so the center of P does not lie in the interior of Q); it follows that the lengths of the edges and diagonals in Q are all smaller than d. Thus d must the be the base of the isosceles triangle in Q, from which it follows that the isosceles triangle is $\Delta P_a P_{a+k/2} P_{a+k}$, and so 2|k. Repeating this process on the legs of isosceles triangle $(\overline{P_a P_{a+k/2}}, \overline{P_{a+k} P_{a+k/2}})$, it follows that $k=2^m$ for some positive integer m (if we allow degeneracy, then we can also let m=0). Now take the isosceles triangle $P_x P_y P_z$, $0 \leq x < y < z < n$ in the triangulation that contains the center of \mathcal{P} in its interior; if a diagonal passes through the center, select either of the isosceles triangles with that diagonal as an edge. Without loss of generality, suppose $P_x P_y = P_y P_z$. From our previous result, it follows that there are 2^a edges of P on the minor arcs of $P_x P_y$, $P_y P_z$ and 2^b edges of P on the minor arc of $P_z P_x$, for positive integers a, b. Therefore, we can write

$$n = 2 \cdot 2^a + 2^b = 2^{a+1} + 2^b$$

so n must be the sum of two powers of 2. We now claim that this condition is sufficient. Suppose without loss of generality that $a + 1 \ge b$; then we rewrite this as

$$n = 2^b(2^{a-b+1} + 1).$$

Lemma 1: All regular polygons with $n = 2^k + 1$ or n = 4 have triangulations that meet the conditions. By induction, it follows that we can cover all the desired n. For n =3, 4, this is trivial. For k>1, we construct the diagonals of equal length $\overline{P_0P_{2^{k-1}}}$ and $P_{2^{k-1}+1}P_0$. This partitions \mathcal{P} into 3 regions: an isosceles $\triangle P_0P_{2^{k-1}}P_{2^{k-1}+1}$, and two other regions. For these two regions, we can recursively construct the isosceles triangles defined above in the second paragraph. It follows that we have constructed $2(2^{k-1}-1)$ + $(1) = 2^k - 1 = n - 2$ isosceles triangles with non-intersecting diagonals, as desired. Lemma 2: If a regular polygon with n sides has a working triangulation, then the regular polygon with 2n sides also has a triangulation that meets the conditions. We construct the diagonals $\overline{P_0P_2}$, $\overline{P_2P_4}$, ... $\overline{P_{2n-2}P_0}$. This partitions \mathcal{P} into n isosceles triangles of the form $\triangle P_{2k}P_{2k+1}P_{2k+2}$, as well as a central regular polygon with n sides. However, we know that there exists a triangulation for the n-sided polygon that yields n-2 isosceles triangles. Thus, we have created (n)+(n-2)=2n-2 isosceles triangles with non-intersecting diagonals, as desired. In summary, the answer is all n that can be written in the form $2^{a+1} + 2^b$, $a, b \ge 0$. Alternatively, this condition can be expressed as either $n=2^k, k \geq 2$ (this is the case when a+1=b) or n is the sum of two distinct powers of 2, where $1=2^0$ is considered a power of 2.

[Source]:

USAMO 2008

</example math problem>

Summarization

The problem asks for the possible values of n for a regular n-sided polygon that can be completely triangulated into isosceles triangles using non-intersecting diagonals. The solution involves analyzing the properties of the diagonals forming isosceles triangles and deducing that n can be expressed in terms of powers of 2.

Difficulty

Rating: 7

Reason: The problem involves understanding properties of isosceles triangles in the context of polygon triangulation and requires critical reasoning to establish relationships between the number of sides and powers of 2, making it more complex than typical undergraduate-level problems.

```
1404
           === report over ===
1405
1406
1407
           <math problem>
           [QUESTION]:
1408
           ***Question***
1409
           [SOLUTION]:
1410
           ***Solution***
1411
           [SOURCE]:
1412
           ***SOURCE***
1413
           </math problem>
1414
1415
1416
           Your answer should be in JSON format for example:
           "json
1417
1418
           "Rating": YOUR RATING,
1419
           "Reason": YOUR JUSTIFICATION,
1420
1421
1422
1423
```

D.7 TOPIC DOMAIN

Meta Prompt for Topic Domain Annotation

I am a teacher, and I have some high-level Olympiad math problems. I want to categorize the domain of these math problems.

OBJECTIVE

- 1. Summarize the math problem in a brief sentence, describing the concepts involved in the math problem.
- 2. Categorize the math problem into specific mathematical domains. Please provide a classification chain, for example, Applied Mathematics -> Probability -> Combinations.

The following is a basic classification framework in the field of mathematics.

<math domains>

1424

1425 1426

1427 1428

1429 1430

1431

1432

1433

1434

1435

1436

1437 1438

1439

1441

1444

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

</math domains>

STYLE

1440 Data report.

TONE

Professional, scientific.

1442 Professional, scie # AUDIENCE #

Students. Enable them to better understand the domain of the problems.

RESPONSE: MARKDOWN REPORT

Summarization

[Summarize the math problem in a brief paragraph.]

Math domains

[Categorize the math problem into specific mathematical domains, including major domains and subdomains.]

ATTENTION

- The math problem can be categorized into multiple domains, but no more than three. Separate the classification chains with semicolons(;).
- Your classification MUST fall under one of the aforementioned subfields; if it really does not fit, please add "Other" to the corresponding branch. For example: Algebra -> Intermediate Algebra -> Other. Only the LAST NODE is allowed to be "Other"; the preceding nodes must strictly conform to the existing framework.
- The math domain must conform to a format of classification chain, like "Applied Mathematics -> Probability -> Combinations".

```
1458
             - Add "=== report over ===" at the end of the report.
1459
1460
              <example math problem>.
1461
              [Question]
              Determine the greatest real number C, such that for every positive integer n \geq 2, there exists
1462
              x_1, x_2, ..., x_n \in [-1, 1], so that
1463
                                                \prod_{1 \le i < j \le n} (x_i - x_j) \ge C^{\frac{n(n-1)}{2}}
1464
1465
1466
1467
              </example math problem>
1468
             ## Summarization
             The problem seeks to find the greatest real number C such that, for every integer n > 2, there
1469
              exist real numbers x_1, x_2, \ldots, x_n \in [-1, 1] satisfying the inequality \prod_{1 \le i \le j \le n} (x_i - x_j) \ge 1
1470
             C^{rac{n(n-1)}{2}}. This involves maximizing C to ensure the product of all pairwise differences
1471
1472
              among n points in the interval [-1, 1] is at least C raised to the power of the number of such
1473
             pairs, \frac{n(n-1)}{2}.
1474
             ## Math domains
              Algebra -> Intermediate Algebra -> Inequalities; Discrete Mathematics -> Combinatorics
1476
              === report over ===
1477
              <example math problem>
1478
              [Question]
1479
              Given integer n \geq 2. Find the minimum value of \lambda, satisfy that for any real numbers a_1, a_2,
1480
              \cdots, a_n and b,
1481
                                       \lambda \sum_{i=1}^{n} \sqrt{|a_i - b|} + \sqrt{n \left| \sum_{i=1}^{n} a_i \right|} \geqslant \sum_{i=1}^{n} \sqrt{|a_i|}.
1482
1483
1484
              </example math problem>
1485
              ## Summarization
1486
              Let n \geq 2 be an integer. The problem seeks the minimum value of \lambda such that for any
1487
             real numbers a_1, a_2, \ldots, a_n and b, the inequality \lambda \sum_{i=1}^n \sqrt{|a_i - b|} + \sqrt{n |\sum_{i=1}^n a_i|} \geqslant
1488
              \sum_{i=1}^n \sqrt{|a_i|} holds. The goal is to find the smallest \lambda that ensures this inequality is satisfied
1489
              for all possible choices of a_i and b.
1490
              ## Math domains
1491
              Algebra -> Intermediate Algebra -> Inequalities;
1492
             Calculus -> Differential Calculus -> Applications of Derivatives.
1493
              === report over ===
1494
1495
1496
              <math problem>
1497
              [Ouestion]
              ***Question***
1498
              [Solution]
1499
              ***Solution***
1500
              [Source]
1501
              ***Source***
1502
              </math problem>
1503
              Your answer should be in JSON format for example:
              """ json
1507
              "Summary": "YOUR SUMMARY",
              "Domains": [domain1, ...]
1509
1510
              }
"""
1511
```

D.8 THEOREM VARIANT CONSTRUCTION

Meta Prompt for Constructing Theorem Variants

I'm going to give you a math proof question and its solution. Your task is to follow the steps below to write a new question based on the given one. Here is the original question and solution:

"'question > question>
"',
"'solution < solution>

Please follow these steps:

- 1. The original question asks to "prove or disprove" a statement, where the statement can be "proved". Please write a new question by negating the original statement, so that it can now be "disproved". For example, if the original statement is x=y, you may change it to x< y or $x \neq y$; if the statement is "there exists xxx", you may change it to "there does not exist xxx". When negating the original question, you should make minimal changes, i.e. leave as much background information unchanged as possible.
- 2. After changing the question, the solution should be changed accordingly. You do not have to write a new solution, and the original solution can probably be reused. For example, if the original question asks to prove x=y and the new question asks to prove x< y, you may simply add a step to the original proof like "since we proved x=y, the statement x< y is disproved". However, check the wording of the solution so that it tries to "prove" the statement at first, and then naturally transit to finding that it cannot be proved, but can be disproved instead.
- 3. The original solution ends with " $boxed{proved}$ ". Your new solution should end with " $boxed{disproved}$ ".

Output the new question and solution in two blocks:

""question
new question
""
""solution
corresponding solution
""

D.9 SYSTEM PROMPT FOR RL TRAINING

System Prompt for RL Training

D.10 PROCESS EVALUATION FRAMEWORK

Meta Prompt for Process Evaluation

You are an expert in scoring solutions for mathematical proof questions. The following question asks to prove or disprove a statement, where the statement may be either true or false. The test subject is asked to end their proof with \boxed{proved} if they prove the statement to be true, and \boxed{disproved} if they prove the statement to be false.

The question: "'<question>"

The ground truth of the statement:

"'<answer>"

The test subject's solution:

"" <solution>""

Your task is to evaluate the proof's quality and assign a score from 0 to 1 based on four criteria: logical validity (40%), completeness (30%), correctness (20%), and clarity (10%). Instructions:

- 1. Analyze the proof step by step.
- 2. For each criterion:
- Logical Validity: Check if each step follows logically from the previous one. Flag any logical errors.
- Completeness: Verify if all necessary cases and steps are included to prove the theorem.
- Correctness: Confirm if the final conclusion is correct.
- Clarity: Assess if the proof is clear, unambiguous, and well-explained.
- 3. Assign a sub-score (0 to 1) for each criterion and compute the total score using the weights: $(0.4 \times \text{validity}) + (0.3 \times \text{completeness}) + (0.2 \times \text{correctness}) + (0.1 \times \text{clarity})$.
- 4. Provide a brief explanation (2-3 sentences) summarizing any errors or issues and justifying the score.

Final output format:

```
{
"score": float,
"validity": float,
"completeness": float,
"correctness": float,
"clarity": float,
"explanation": str
}
```

where "score" is the total score, and "validity", "completeness", "correctness", "clarity" are the subscores.