

Reliable Use of Lemmas via Eligibility Reasoning and Section-Aware Reinforcement Learning

Zhikun Xu^{1*}, Xiaodong Yu², Ben Zhou¹, Jiang Liu², Jialian Wu²,
Ze Wang², Ximeng Sun², Hao Chen², Zicheng Liu²

¹Arizona State University, ²Advanced Micro Devices, Inc.
zhikunxu@asu.edu

Abstract

Recent large language models (LLMs) perform strongly on mathematical benchmarks yet still misapply lemmas by importing conclusions whose assumptions are not satisfied. We study this failure mode through lemma-judging: given a statement and a candidate lemma, the model predicts (i) whether the lemma’s preconditions hold and (ii) whether its conclusion is useful, from which the final usefulness label is derived. We present RULES, a two-section output protocol trained with reinforcement learning and section-aware loss masking so that negative credit is assigned to the section responsible for an error when fine-grained blame is available. Training uses diverse natural-language and formal proof corpora, and robustness is evaluated on a held-out perturbation suite that minimally edits statements to break lemma applicability while preserving coherence. Across four base models, RULES improves in-domain lemma-judging and achieves its largest gains on this applicability-breaking test, consistent with better gating against inapplicable lemmas. Transfer to end-to-end benchmarks is mixed but broadly competitive with a single-label RL baseline, reflecting a precision/recall tradeoff rather than a claim of end-to-end safety.

1 Introduction

Recent large language models (LLMs) achieve strong scores on mathematical benchmarks and can often generate proof-like text (Welleck et al., 2022; OpenAI et al., 2024; Yang et al., 2024; Shao et al., 2024; Li et al., 2025b; Guo et al., 2025). Yet their behavior frequently reflects pattern imitation rather than principled reasoning: models tend to reproduce familiar solution paths, and they can import a lemma’s conclusion without validating its hypotheses. Such behavior exploits surface corre-

lations in the training distribution and undermines trustworthiness in mathematical reasoning.

Mathematicians typically verify a lemma’s hypotheses before invoking it in a proof (Frieder et al., 2024). We operationalize this practice as a structured prediction problem (lemma-judging): given a statement-lemma pair, the model must output: ① a precondition check (validating if the lemma’s assumptions hold in the current context) and ② a conclusion-utility check (assessing if the lemma is actually helpful for the proof goal), from which a useful decision is derived. Our goal is not to claim a comprehensive, end-to-end safety guarantee, but to target a specific failure mode in mathematical reasoning: using a lemma whose hypotheses do not hold. We hypothesize that explicit eligibility reasoning with aligned credit assignment will reduce reliance on spurious surface cues and improve robustness to applicability-breaking perturbations. In this paper, we train models only on the lemma-judging task and then evaluate the resulting models directly on end-to-end math benchmarks to test transfer. Operationally, RULES can be viewed either as a training-time regularizer for mathematical reasoning or as a verifier/reranker that screens candidate lemmas in a larger pipeline.

We propose a training framework, RULES (Reliable Use of Lemmas via Eligibility Reasoning and Section-Aware Reinforcement Learning). The training is based on *Group Relative Policy Optimization* (GRPO) (Shao et al., 2024) with a section-aware loss masking mechanism, which assigns penalty to the precise section responsible for failure, rather than the entire rollout. This design aligns the supervision with the cognitive decomposition the task demands, encouraging models to learn the check-before-use discipline rather than overfitting to label-level shortcuts.

Empirical results over various models and diverse evaluation datasets have shown that RULES improves in-domain lemma-judging over a vanilla

*Work done during research internship at AMD.

model and a single-label RL baseline. The margin is larger on the held-out applicability-breaking perturbation suite with at most 23 points, indicating better gating of lemma use. On end-to-end benchmarks, average performance is broadly better or on par with the single-label RL baseline. Ablations show that both the two-section outputs and section-aware RL are necessary: entangling fine-grained supervision on checks or excluding training-time perturbations supervision substantially reduces the robustness of eligibility reasoning.

2 Related Works

Many works have explored pre-training on math corpora coupled with chain-of-thought prompting, which equips LLMs with domain-specific understanding and step-by-step reasoning (Wang et al., 2025a). Reinforcement learning advances, such as GRPO, have offered significant gains on benchmarks like GSM8K and MATH (Shao et al., 2024). However, brittleness to premise order and problem phrasing remains a challenge (Chen et al., 2024; Huang et al., 2025). Self-correction methods have also improved LLM mathematical reasoning by enabling self-refinement (Xiong et al., 2025; Yan et al., 2025). Benchmarks such as the Open Proof Corpus (Dekoninck et al., 2025) and MathArena (Balunovic et al., 2025) reveal a gap between the accuracy of final answers and the validity of proofs. Knowledge-graph agents have been explored to improve success rates (Li et al., 2025a), while LeanDojo (Yang et al., 2023) and other works (Mikuła et al., 2023; Tao et al., 2025) enhance premise selection. Recent innovations (Mukherjee et al., 2025; Pan et al., 2025; Wang et al., 2025b) further align supervision with intermediate reasoning, improving error detection and reflective reasoning. Our work builds on existing insights and focuses on structured intermediate reasoning that aims to improve generalizable LLM reasoning on math.

3 Methodology: RULES

RULES has two components: ① a **two-section lemma-judging task** that requires explicit eligibility checks before use, and ② **section-aware reinforcement learning** that aligns credit assignment with these checks.

3.1 Lemma-judging Task

Given an input $x = (L, S)$ with lemma L and statement S , the model produces a structured output with two sections $y = (r_{pre}, l_{pre}) \oplus (r_{con}, l_{con})$ where r_{pre} and r_{con} are free-form rationales for the precondition and conclusion-utility checks, and the section (predicted) labels, i.e., $l_{pre}, l_{con} \in \{0, 1\}$, are extracted judgments. The final usefulness decision is a deterministic aggregation of the precondition and conclusion checks $l_{use} = l_{pre} \wedge l_{con} \in \{0, 1\}$. This task specification makes the precondition and conclusion checks explicit and auditable, and a minimal schema with section sentinels ensures unambiguous parsing, which is illustrated in Appendix A by prompts.

3.2 Section-Aware Reinforcement Learning

The training is adapted from GRPO. Specifically, for each input x , given a rollout as $y^i = y_{pre}^i \oplus y_{con}^i = (r_{pre}^i, l_{pre}^i) \oplus (r_{con}^i, l_{con}^i)$, if the precondition check is correct but the aggregated final label is incorrect, which implicitly means that the conclusion check must be wrong, we will only keep the conclusion check section, i.e., setting the loss mask of (r_{con}^i, l_{con}^i) to be 1 and others as 0. For the RL update, the gradient function in each group G is illustrated as: $\nabla_{\theta} L_{section-aware} = -\frac{1}{|G|} \sum_{i=1}^{|G|} \frac{1}{|y^i|} \sum_{j=1}^{|y^i|} m_j^i \tilde{A}_j^i \nabla_{\theta} \log \pi_{\theta}(y_j^i | x, y_{<j}^i)$, where m_j^i is the j -th token mask for the i -th rollout and other notations in this equation follow the default GRPO setting such as estimated advantages \tilde{A}_j^i . By this, we are penalizing the incorrect sections only and supervising models in more fine-grained signals on certain wrong tokens, enforcing "check-before-use" as a safety filter for theorem application. Crucially, this objective-level modification operates independently of the underlying model architecture or scale. The computational overhead is restricted to negligible token-masking operations during the loss calculation, making RULES compatible with any LLM backbones.

4 Experimental Settings

4.1 Models

We have mainly experimented with four LLMs: **DeepMath-1.5B** (He et al., 2025), **OLMo-2-1124-7B-Instruct** (OLMo et al., 2025), **Qwen2.5-Math-7B-Instruct** (Yang et al., 2024), and **Llama-3.1-8B-Instruct** (Grattafiori et al., 2024). For each model, we report three variants: ① **Vanilla**: di-

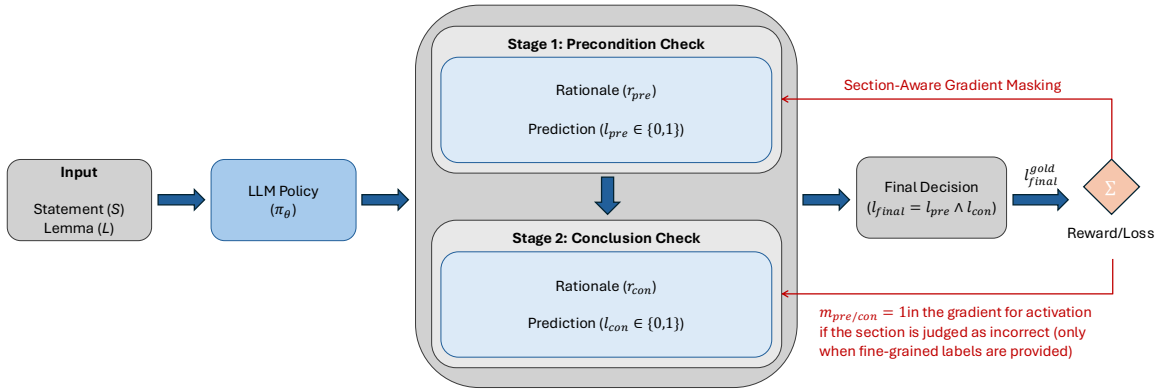


Figure 1: Overview of RULES Training Framework. As mentioned in Section 3.2, if the fine-grained label (precondition/conclusion check) is provided, we could know which section of tokens are to blame with negative rewarding.

Models \ Tests		In-Domain Test				Perturbation Test			E2E Benchmark			
		NP	NLPS	IS	IL	DT	PA	CM	IM	MP-hard	MP-simple	TQ-Math
DeepMath	vanilla	60.8	51.5	51.6	31.9	24	39.1	45.7	76	53	72	55.4
	GRPO	77.5	77.2	67.6	65.8	35.6	37.9	63.1	73	54.5	72.8	56.1
	RULES	88.5	86.2	75.8	89.5	47	39.5	46	71	56.3	72.8	57.7
OLMO2	vanilla	77.6	69.1	53.9	64.9	47.6	8.4	52.3	29	8.6	18.6	27.1
	GRPO	85.2	84.2	72.2	83.6	46.8	7.8	48.1	23	9.3	21.9	30.3
	RULES	84.9	83.4	76.6	82.5	60.5	6.7	48.2	24	10.4	23.3	26.9
Qwen2.5	vanilla	63.9	71.9	55.5	55.8	71.9	32.4	59.5	55	33.7	59.5	49.8
	GRPO	76.5	83.7	74.2	80.4	81.2	35.6	60.5	54	33.7	59.5	50
	RULES	88.2	85.8	81.3	92.7	87.8	32	61.3	53	35.1	62.4	50.5
Llama-3.1	vanilla	61.3	75.1	54.7	52.3	89.2	11.5	56.5	22	13.3	38.7	33.5
	GRPO	88.6	86.5	85.9	89.8	80.9	13.2	50.7	18	13.3	36.2	32.6
	RULES	93.2	86.4	86.3	93.6	90.9	12.8	49.3	16	12.9	36.9	33

Table 1: **Main Results.** For in-domain test, NP, NLPS, IS, and IL stands for NaturalProofs, NLPS, Premise Selection in Isabelle, IMO-Lemma derived test sets. DT stands for DeepTheorem-sampled perturbation test set. PA, CM, IM, MP, and TQ stand for Putnam-Axiom, CounterMATH, IneqMath, MATH-Perturb, and TheoremQA. MP-hard and MP-simple stand for hard and simple perturbations based on the paper. TQ-Math is the math subset of the TQ benchmark. The CounterMATH is using F1 due to imbalanced data distribution mentioned in Li et al. (2025b) and the rest of data is using exact-match based accuracy.

rect inference with our two-section schema. ② **GRPO**: default GRPO optimization with only the final usefulness label (single-label objective) using the vanilla prompt. ③ **RULES**: GRPO with section-aware masking by two-section schema.

4.2 Training Settings

The training data are from natural language and formal proof corpora, i.e., **NaturalProofs** (Welleck et al., 2021), **NLPS** (Natural Language Premise Selection) (Ferreira and Freitas, 2020), **Premise Selection in Isabelle** (Mikuła et al., 2023), and **IMO-lemmas** (Liang et al., 2025), pairing statement with candidate lemmas and binary labels showing the usefulness of lemmas. We additionally construct a perturbation subset of 744 applicability-breaking examples after filtering. Because only this perturba-

tion subset provides reliable section-local supervision, section-aware masking is intentionally sparse, at roughly 5.7% of the core corpus scale. The remaining updates rely on the standard rollout-level binary reward. More training settings and designs are in Appendix B.

Perturbation Design. As our goal is to separate surface similarity from true eligibility, we also add some perturbation data in training by making a lemma inapplicable while keeping the statement minimally changed and mathematically coherent. The reason of only providing labels for the precondition check when constructing perturbations is two fold: ① Precondition satisfaction is directly falsifiable from the statement-lemma pair and independent of proof strategy, yielding reliable supervi-

Ablations	In-Domain Test				Perturbation Test				E2E Benchmark			
	NP	NLPS	IS	IL	DT	PA	CM	IM	MP-hard	MP-simple	TQ-Math	
vanilla	63.9	71.9	55.5	55.8	71.9	32.4	59.5	55	33.7	59.5	49.8	
GRPO	76.5	83.7	74.2	80.4	81.2	35.6	60.5	54	33.7	59.5	50	
two-section-onetime	87	85.4	84.4	92.7	73.1	32.4	60.5	49	35.5	60.9	49.5	
w/o perturbation data	93.2	87.4	82.4	94.4	19.1	32.2	56.4	49	34.4	59.1	50.5	
de-noised perturbation data	86	83.9	83.6	95	74.3	34.1	59.6	54	35.8	58.8	50.9	
RULES	88.2	85.8	81.3	92.7	87.8	32	61.3	53	35.1	62.4	50.5	

Table 2: **Ablation Results.** *two-section-onetime* refers to output two checks predictions at last. *w/o perturbation data* and *de-noised perturbation data* refers to removing perturbation data in training and filter the unqualified perturbation training data by the o3 model and human validation, respectively. The experiments are done with Qwen2.5-Math-7B-Instruct.

sion with low ambiguity. ② Conclusion utility is inherently strategy-dependent (e.g., proof by counterexample), making it costly and noisy to annotate at scale.

Reward Design. For a fair comparison with baselines, the reward is rule-based and binary, i.e., 1 for correct and -1 for incorrect. Moreover, outputs that fail to parse under the two-section schema receive a format penalty of -2.

4.3 Evaluation Settings

We evaluate three axes that mirror our intended use: in-domain lemma-judging, robustness to applicability-breaking perturbations, and transfer to end-to-end mathematical problem solving. All results are calculated by sampling 7 outputs and majority-voting the extracted answers by self-consistency, which is noted as **SC@7**. More evaluation details are shown in Appendix C.

In-domain Evaluation We use held-out splits drawn from the training data sources, which are named as **naturalproofs-test**, **nlps-test**, **isabelle-test**, and **IMO-lemma-test**.

Perturbation Test Similar to the goal and process of training perturbation data curation, we also evaluate on an out-of-domain perturbation suite constructed by minimally editing statements to violate lemma preconditions. But the raw data is sampled from DeepTheorem (Zhang et al., 2025), which is noted as **deeptheorem-perturbation**. More details about perturbation data quality are in Appendix D.

End-to-end Benchmarks Besides the task-related evaluations, we also assess transfer on diverse mathematical tasks, including competition-style problems (**Putnam-Axiom** (Gulati et al., 2024)), inequality reasoning (**IneqMath** (Sheng

et al., 2025b)), counterexample-driven questions (**CounterMATH** (Li et al., 2025b)), perturbation-aligned suites (**MATH-Perturb** (Huang et al., 2025)), and theorem-centric math QA (**TheoremQA** (Chen et al., 2023)).

5 Analysis

5.1 Main Results

In-domain lemma-judging. Across all four base models, RULES improves in-domain lemma-judging over both *vanilla* and the single-label *GRPO* baseline. Gains are consistent on most splits and especially pronounced on the IMO-derived test, with OLMO2 showing comparable performance to GRPO on three splits and clear improvements on the Isabelle split. These results indicate that enforcing explicit precondition and conclusion checks could help models improve understanding of lemma applicability.

Robustness to applicability-breaking perturbations. On the out-of-domain perturbation suite, RULES consistently outperforms *vanilla* and *GRPO* across all four models, with sizable margins for 23+ points at most. This pattern supports the intended effect of training-time perturbation data and section-aware masking: the model more reliably gates off familiar but inapplicable lemmas under minimal statement edits that invalidate the preconditions.

Mixed E2E transfer. On downstream tasks, RULES performs broadly on par with single-label *GRPO*, yielding gains on perturbation-aligned subsets while showing occasional regressions on competition-style and inequality-based problems. We interpret these regressions as a **performance-rigorosity tradeoff**: whereas standard models may sometimes “get lucky” by hallucinating applicability to force a solution, RULES enforces strict

precondition checks, which is illustrated in Appendix G. This discipline prioritizes precision and safety against invalid reasoning, naturally lowering recall on tasks where looser constraints might inadvertently permit a correct answer. Consequently, gains are most pronounced on perturbation-aligned suites where robustness against invalid inputs is paramount. Conversely, constructive tasks like CounterMATH require example synthesis, a capability distinct from eligibility gating. Thus, RULES functions as a necessary robustness layer, ensuring models do not attempt proofs with fundamentally inapplicable premises.

5.2 Ablations

Two-time v.s. one-time prediction. Replacing the sequential two-section protocol with a “two-section-onetime” variant (emitting both judgements only at the end) yields strong in-domain accuracy but notably weaker robustness on the perturbation suite. This indicates that the staged structure and section-aware alignment matters for learning to withhold lemmas when preconditions are not satisfied.

Role of training-time perturbation data. Removing perturbation data during training preserves in-domain scores but collapses robustness on the out-of-domain suite, confirming that curated negatives are essential to teach the precondition gate. However, we emphasize that this robustness is not merely a function of data scaling. The architecture contributes significantly: RULES improves over the data-equivalent “two-section-onetime” baseline by 14.7 points on the perturbation test (Table 2). This substantial margin confirms that while the data provides the necessary signals, the **Section-Aware RL mechanism** is indispensable for effectively leveraging them to prevent superficial pattern matching.

Effect of "de-noised" perturbation data. Filtering low-quality perturbations leads to mixed outcome: modest improvements on some end-to-end metrics and in-domain tests, but still below RULES on the perturbation suite. This evidence thus far does not support that "clean fine-grained labels" always improve overall math reasoning; rather, there appears to be a trade-off between strict cleaning and robustness coverage.

6 Conclusion

We propose RULES, aligning lemma use with an explicit check-before-use protocol and section-

aware RL. RULES yields consistent in-domain gains, large robustness improvements under precondition-breaking perturbations, and parity or small gains on end-to-end tasks, which is best understood as a precision/recall tradeoff rather than a general safety improvement. Ablations show the staged two-section output and section-aligned credit assignment are both required.

Limitations

Task scope and supervision. Our study focuses on *lemma-judging* rather than full proof generation. The supervision signals are tailored to the specific task, which may bias models toward unexpected behaviors.

Label coverage and noise. In the perturbation data, only the precondition label is available while the conclusion-utility label is omitted to avoid strategy dependence. This asymmetric supervision can introduce bias and limit learning of usefulness beyond eligibility. In addition, rewards are binary with a format penalty, which simplifies optimization but may underutilize graded signals.

Model and data coverage. Experiments cover four mid-sized models and a specific suite of math corpora (NaturalProofs, NLPS, Isabelle premise selection, and IMO-lemma). This leaves gaps across subfields of mathematics and automated theorem proving techniques, and languages other than English. While the eligibility reasoning mechanism is theoretically language-agnostic, the extension to multilingual settings or other formal proof assistants (e.g., Lean, Coq) remains a direction for future work.

References

- Mislav Balunovic, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. 2025. [Matharena: Evaluating LLMs on uncontaminated math competitions](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023. [TheoremQA: A theorem-driven question answering dataset](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7889–7901, Singapore. Association for Computational Linguistics.
- Xinyun Chen, Ryan A Chi, Xuezhi Wang, and Denny Zhou. 2024. [Premise order matters in reasoning with](#)

- large language models. In *Proceedings of the 41st International Conference on Machine Learning*, pages 6596–6620.
- Jasper Dekoninck, Ivo Petrov, Kristian Minchev, Miroslav Marinov, Maria Drencheva, Lyuba Konova, Milen Milenov Shumanov, Kaloyan Tsvetkov, Nikolay Drenchev, Lazar D. Todorov, Kalina Nikolova, Nikolay Georgiev, Vanesa Kalinkova, Margulan Ismoldayev, Mislav Balunovic, and Martin Vechev. 2025. [The open proof corpus: A large-scale study of LLM-generated mathematical proofs](#). In *2nd AI for Math Workshop @ ICML 2025*.
- Deborah Ferreira and André Freitas. 2020. [Natural language premise selection: Finding supporting statements for mathematical text](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2175–2182, Marseille, France. European Language Resources Association.
- Simon Frieder, Julius Berner, Philipp Petersen, and Thomas Lukasiewicz. 2024. [Large language models for mathematicians](#). *Preprint*, arXiv:2312.04556.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Aryan Gulati, Brando Miranda, Eric Chen, Emily Xia, Kai Frønsdal, Bruno de Moraes Dumont, and Sanmi Koyejo. 2024. [Putnam-AXIOM: A functional and static benchmark for measuring higher level mathematical reasoning](#). In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*.
- Dadi Guo, Jiayu Liu, Zhiyuan Fan, Zhitao He, Haoran Li, Yumeng Wang, and Yi R. Fung. 2025. [Mathematical proof as a litmus test: Revealing failure modes of advanced large reasoning models](#). *Preprint*, arXiv:2506.17114.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. [Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning](#). *Preprint*, arXiv:2504.11456.
- Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun Chen, Chiyuan Zhang, and Mengdi Wang. 2025. [MATH-perturb: Benchmarking LLMs' math reasoning abilities against hard perturbations](#). In *Forty-second International Conference on Machine Learning*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#). In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Vincent Li, Tim Knappe, Yule Fu, Kevin Han, and Kevin Zhu. 2025a. [Automating mathematical proof generation using large language model agents and knowledge graphs](#). *arXiv preprint arXiv:2503.11657*.
- Yinghui Li, Jiayi Kuang, Haojing Huang, Zhikun Xu, Xinnian Liang, Yi Yu, Wenlian Lu, Yangning Li, Xiaoyu Tan, Chao Qu, Ying Shen, Hai-Tao Zheng, and Philip S. Yu. 2025b. [One example shown, many concepts known! counterexample-driven conceptual reasoning in mathematical LLMs](#). In *Forty-second International Conference on Machine Learning*.
- Zhenwen Liang, Linfeng Song, Yang Li, Tao Yang, Feng Zhang, Haitao Mi, and Dong Yu. 2025. [Towards solving more challenging IMO problems via decoupled reasoning and proving](#). *arXiv preprint*.
- Maciej Mikula, Szymon Antoniak, Szymon Tworkowski, Bartosz Piotrowski, Albert Jiang, Jin Peng Zhou, Christian Szegedy, Łukasz Kuciński, Piotr Miłoś, and Yuhuai Wu. 2023. [Magnushammer: A transformer-based approach to premise selection](#). In *The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS'23*.
- Sagnik Mukherjee, Abhinav Chinta, Takyoun Kim, Tarun Anoop Sharma, and Dilek Hakkani Tur. 2025. [Premise-augmented reasoning chains improve error identification in math reasoning with llms](#). In *Forty-second International Conference on Machine Learning*.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, and 21 others. 2025. [2 olmo 2 furious](#). *Preprint*, arXiv:2501.00656.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Zhuoshi Pan, Yu Li, Honglin Lin, Qizhi Pei, Zinan Tang, Wei Wu, Chenlin Ming, H. Vicky Zhao, Conghui He, and Lijun Wu. 2025. [LEMMA: Learning from errors for MatheMatical advancement in LLMs](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 11615–11639, Vienna, Austria. Association for Computational Linguistics.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan

- Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *Preprint*, arXiv:2402.03300.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025a. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference on Computer Systems, EuroSys '25*, page 1279–1297, New York, NY, USA. Association for Computing Machinery.
- Jiayi Sheng, Luna Lyu, Jikai Jin, Tony Xia, Alex Gu, James Zou, and Pan Lu. 2025b. Solving inequality proofs with large language models. In *The Thirtieth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yicheng Tao, Haotian Liu, Shanwen Wang, and Hongteng Xu. 2025. Learning an effective premise retrieval model for efficient mathematical formalization. In *2nd AI for Math Workshop @ ICML 2025*.
- Peng-Yuan Wang, Tian-Shuo Liu, Chenyang Wang, Yi-Di Wang, Shu Yan, Cheng-Xing Jia, Xu-Hui Liu, Xin-Wei Chen, Jia-Cheng Xu, Ziniu Li, and 1 others. 2025a. A survey on large language models for mathematical reasoning. *arXiv preprint arXiv:2506.08446*.
- Pinzheng Wang, Juntao Li, Zecheng Tang, Haijia Gui, and 1 others. 2025b. Improving rationality in the reasoning process of language models through self-playing game. In *Forty-second International Conference on Machine Learning*.
- Sean Welleck, Jiacheng Liu, Ronan Le Bras, Hannaneh Hajishirzi, Yejin Choi, and Kyunghyun Cho. 2021. Naturalproofs: Mathematical theorem proving in natural language. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Sean Welleck, Jiacheng Liu, Ximing Lu, Hannaneh Hajishirzi, and Yejin Choi. 2022. Naturalprover: Grounded mathematical proof generation with language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 4913–4927. Curran Associates, Inc.
- Wei Xiong, Hanning Zhang, Chenlu Ye, Lichang Chen, Nan Jiang, and Tong Zhang. 2025. Self-rewarding correction for mathematical reasoning. *arXiv preprint arXiv:2502.19613*.
- Yuchen Yan, Jin Jiang, Yang Liu, Yixin Cao, Xin Xu, Mengdi Zhang, Xunliang Cai, and Jian Shao. 2025. S³cmath: Spontaneous step-level self-correction makes large language models better mathematical reasoners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 25588–25596.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *Preprint*, arXiv:2409.12122.
- Kaiyu Yang, Aidan Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan J Prenger, and Animashree Anandkumar. 2023. Leandojo: Theorem proving with retrieval-augmented language models. *Advances in Neural Information Processing Systems*, 36:21573–21612.
- Ziyin Zhang, Jiahao Xu, Zhiwei He, Tian Liang, Qizhi Liu, Yansi Li, Linfeng Song, Zhenwen Liang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. Deeptheorem: Advancing llm reasoning for theorem proving through natural language and reinforcement learning. *Preprint*, arXiv:2505.23754.

A Prompts

Two-Section Prompt

Lemma: {lemma}
Statement: {statement}

Please analyze step by step whether the lemma is useful for proving or disproving the given statement. You should perform two sequential checks:

Step 1: Precondition Check

First, analyze whether the lemma's preconditions are satisfied in the given statement. After your reasoning, provide your answer within boxed{ } as either boxed{True} (preconditions are satisfied) or boxed{False} (preconditions are not satisfied).

Step 2: Conclusion Check

Next, analyze whether the lemma's conclusion is helpful for proving or disproving the given statement. After your reasoning, provide your answer within boxed{ } as either boxed{True} (conclusion is helpful) or boxed{False} (conclusion is not helpful).

Please use the same section titles as above (i.e., **Step 1: Precondition Check** and **Step 2: Conclusion Check**) to start the two checks respectively and STRICTLY follow the following format:

- First, provide your reasoning for the precondition check
- Then output boxed{True} or boxed{False} for the precondition check
- Then provide your reasoning for the conclusion check
- Finally output boxed{True} or boxed{False} for the conclusion check

Vanilla Prompt

Lemma: {lemma}
Statement: {statement}

Please reason step by step about whether the above lemma is useful for proving or disproving its following statement, and then put your final answer (i.e., True or False) within boxed.

Two-Section Onetime Prompt

Lemma: {lemma}
Statement: {statement}

Please think step by step whether the lemma is useful for proving or disproving the given statement. Specifically, please 1) check whether the lemma's preconditions are satisfied in the given statement, and 2) check whether the lemma's conclusion is helpful for proving the given statement. Then put your final answers for precondition check and conclusion check, i.e., True or False for each check, within boxed{} (For example, boxed{True, False} means the lemma's preconditions are satisfied but its conclusion is not helpful for proving the given statement).

Precondition Perturbation Prompt for o3 (DeepTheorem)

Based on the given statement and its proof, could you help me perturb the statement with the minimal edits so that one certain precondition of some lemma used in the proof is no longer satisfied? If the perturbation is not possible, please output suitability: False. Otherwise, return suitability: True, and then extract the affected lemma in the proof, generate the perturbed statement, and the rationale of the perturbation.

Critical Requirements:

1. **MATHEMATICAL COHERENCE:** The perturbed statement must be mathematically well-defined and meaningful. Do not create statements that are conceptually impossible or nonsensical (e.g., "subgroups of semigroups").

2. **GENUINE INAPPLICABILITY:** The chosen lemma must be fundamentally unusable in ANY reasonable proof approach for the perturbed statement. Avoid superficial changes where:

- The lemma applies to a subset/restriction of the new domain
- The lemma can be used indirectly through standard correspondences
- The statement can be reduced back to the lemma's domain

3. **VERIFICATION TEST:** Before finalizing, ask: "Could the chosen lemma still contribute meaningfully to proving or disproving the perturbed statement through any mathematical pathway?" If yes, the perturbation is ineffective.

Rules:

- Copy the OriginalStatement verbatim except for that one span.
- Do NOT strengthen assumptions, introduce unrelated concepts/symbols, or make changes that render the statement trivial.
- Prefer the simplest change that creates a fundamental barrier to lemma applicability.
- The perturbation should target the lemma's core assumptions or domain in a way that cannot be circumvented.

Ineffective Perturbation Patterns to Avoid:

- Domain changes where standard correspondences exist (e.g., \mathbb{R} exponential \rightarrow \mathbb{C} exponential when differentiating w.r.t. real variables)
- Structure downgrades that still allow the target concept (e.g., group \rightarrow semigroup but keeping "subgroup" terminology)
- Changes that only affect surface terminology but not mathematical substance

Return ONLY this schema:

Suitability: {bool}
AffectedLemma: {lemma-content}
OriginalStatement: {original statement}
PerturbedStatement: {perturbed statement}
Rationale: {rationale explaining why the lemma becomes fundamentally inapplicable, not just surface-level different}

Now Process:

Statement: {statement}
Proof: {proof}

Precondition Perturbation Validation Prompt for o3

You are evaluating the quality of a mathematical statement perturbation. Your task is to determine if the perturbation effectively breaks the chosen lemma's applicability.

Given Perturbation Data:

Lemma: {lemma}
Original Statement: {statement}
Perturbed Statement: {perturbed statement}
Rationale: {rationale}

Evaluation Criteria (ALL must pass):

1. **MATHEMATICAL COHERENCE CHECK:**

- Is the perturbed statement mathematically well-defined?
- Are all mathematical concepts used consistently and meaningfully?
- Do the mathematical objects and relations make sense together?

Common failures:

- Mixing incompatible concepts (e.g., "subgroups of semigroups")
- Ill-defined operations or structures
- Conceptual contradictions

2. **LEMMA INAPPLICABILITY VERIFICATION:**

Systematically check if the lemma could still be used via:

a) **DIRECT APPLICATION:**

- Can the lemma's hypotheses still be satisfied in the new context?

b) **RESTRICTION/SUBSET APPLICATION:**

- Does the lemma apply to relevant subsets of the new domain?
- Can we restrict to cases where the lemma's conditions hold?

c) **CORRESPONDENCE/ISOMORPHISM:**

- Are there standard mathematical correspondences between the old and new domains?
- Can we map the problem back to where the lemma applies?

d) **INDIRECT/FOUNDATIONAL USE:**

- Could the lemma be used in proving prerequisite results?
- Is the lemma foundational to the theory needed for the perturbed statement?
- Could it be applied to finite quotients, local analysis, etc.?

e) **ANALOGOUS RESULTS:**

- Are there direct analogs of the lemma in the new domain?
- Would the original lemma be a key step in proving such analogs?

3. **FAILURE PATTERN DETECTION:**

Check for these known ineffective patterns:

- **SUPERFICIAL DOMAIN CHANGES:** Changes that look different but preserve mathematical substance
- **TERMINOLOGY SWAPS:** Changing labels without affecting underlying mathematics
- **TRIVIAL GENERALIZATIONS:** Moving to broader categories where restrictions still apply
- **INCOMPLETE DOMAIN SHIFTS:** Changing some but not all relevant aspects of the domain

4. **PROOF PATHWAY ANALYSIS:**

Consider: "In any reasonable attempt to prove or disprove the perturbed statement, could the given lemma contribute meaningfully to the argument?"

If YES to any pathway \rightarrow INEFFECTIVE perturbation

If NO to all pathways \rightarrow EFFECTIVE perturbation

Return your evaluation in this format:

EFFECTIVENESS: [BOOLEAN]
RATIONALE: [Explanation of why the perturbation is effective or ineffective]

Now evaluate the given perturbation:

B More Training Details

Our training experiments are mainly implemented on AMD MI250 and Nvidia H200 GPU servers. The training batch size is 128 and the number of rollout is set to 8. The response length is 8000 and

2048 for reasoning and non-reasoning models respectively. Training checkpoints are used when the percentage of activated section aware loss masking instances of a single batch is close to zero. Our training framework is based on *verl* (Sheng et al., 2025a). The training data distribution is NaturalProofs: 4744, NLPS: 2858, Premise Selection in Isabelle: 4484, and IMO-lemma: 986.

Details of Perturbation Design. Concretely, we target a specific hypothesis of the lemma used in a proof and apply a minimal edit to the statement so that the hypothesis fails ($l_{pre}^{gold} = 0$), without introducing unrelated concepts or trivializing the claim. This is derived from NaturalProofs and processed by o3 model¹. The perturbation data is the only part with the gold precondition label for section-aware masking.

C More Evaluation Details

Our evaluation framework is based on *vllm* (Kwon et al., 2023). The max token length is set to 10k to avoid early truncation and temperature is set to 0.7 for SC@7 across all models. The in-domain and perturbation evaluation data distributions are **naturalproofs-test** (1002), **nlps-test** (676), **isabelle-test** (256), **IMO-lemma-test** (342), and **deeptheorem-perturbation** (839)

D Perturbation Data Quality

To verify the quality of our perturbation suite (where only the precondition label is flipped), we conducted a multi-stage validation. For the training data, we have used O3 model and GPT-OSS-120B to cross-validate them, where both models reached consensus on 84.8% of training data. We only use these 84.8% data as our final perturbation training. Moreover, we authors have also sampled 100 datapoints from it for human verification, and the pass rate is 90%. For perturbation evaluation data, we authors have sampled 100 datapoints after o3-based and human filtering to verify the quality of the final data. Our evaluation guideline follows exactly the principles in the prompt: (i) mathematical coherence (ii) genuine inapplicability (iii) verification for similar pathways by equivalent lemmas. The pass rate for the sampled evaluation data is 95%.

¹<https://platform.openai.com/docs/models/o3>

E Case Study: Preventing Lemma Misapplication

To illustrate how RULES prevents the hallucinations of lemma applicability, we present a case from the Riemannian Geometry domain where a standard vanilla model fails.

Scenario. Statement: Let $G = SU(n)$ and let g be a *left-invariant* metric. Show that any symmetric 2-tensor h satisfying the linearized Einstein equation... must be trivial.

Candidate Lemma: "Since g is *bi-invariant*, it is an Einstein metric... that $Ric_g = \lambda g$ for some $\lambda > 0$."

Analysis. The lemma requires the metric to be *bi-invariant*. However, the problem statement only guarantees that g is *left-invariant*. In differential geometry, left-invariance is a weaker condition than bi-invariance; thus, the lemma's preconditions are not met.

Model Behavior.

- **Baseline (Vanilla):** The model ignores the subtle distinction between "left-invariant" and "bi-invariant," hallucinates that the condition is met, outputs True, and attempts to use the lemma.
- **RULES:** The model correctly identifies the mismatch in the *Precondition Check* section: "*Rationale: The lemma requires bi-invariance, but the statement only assumes left-invariance...*" It outputs `False` for the precondition, effectively gating off the invalid path.

F End-to-end Positive Case: Missing Contraction Precondition

Statement. A non-empty closed convex subset U of a reflexive Banach space with normal structure and a mapping T from U to itself, satisfying $|Tx - Ty| \leq K|x - y|$, then T must have a fixed point.

Rationale. Baseline assumes the Lipschitz constant $K < 1$ without checking — RULES notices K could be 1, making the contraction theorem inapplicable.

Vanilla output. [Truncated] Additionally, U is said to have normal structure. I think normal structure is a property that a Banach space has if every bounded, closed, convex subset that isn't

a singleton has a point that is not diametrically opposed to any other point. In simpler terms, it means that such subsets are not too convex or something like that. This property often relates to the existence of fixed points for certain types of mappings, especially non-expansive ones., the question is whether T must have a fixed point given these conditions. recall that in Banach spaces, there are several fixed point theorems. For non-expansive mappings (which are mappings where $|Tx - Ty| \leq |x - y|$), one of the main theorems is the fixed point theorem by Meir and Keeler, which I think is related to convergence rather than fixed points. the problem here says T is a contraction mapping, since it has a Lipschitz constant $K < 1$ [Truncated]

RULES output. [Truncated] Normal structure in the context of fixed point theory... I think it's a property of a Banach space that ensures that every bounded, closed, convex, and nonempty subset has the fixed point property. Or maybe it's related to the geometric properties of the space? Let me recall. Yeah, in spaces with normal structure, every bounded, closed, convex subset of the space has a nonempty intersection of all its translates, which might imply the existence of a fixed point under certain conditions. Next, the mapping T is given from U to itself, and it satisfies the condition $\|Tx - Ty\| \leq K\|x - y\|$. That sounds like a **Lipschitz** contraction with Lipschitz constant K , right? Wait, but it's not necessarily a contraction mapping because a contraction has $K < 1$. So, does it mean T is a non-expanding mapping with constant K ? But the question is about fixed points. I remember the Banach Fixed Point Theorem which requires a contraction mapping, i.e., $K < 1$. But here K might be greater than or equal to 1, just a Lipschitz constant. [Truncated]

G An Precision/Recall Trade-off Analysis on CounterMath and DeepMath-1.5B

In Table 3, RULES's errors are 5.8x more likely to be over-rejection (saying False when truth is True) than under-rejection (saying True when truth is False). This asymmetry directly reflects the check-before-use training: the model becomes more skeptical of claims, which improves precision (26 false claims correctly caught that baseline missed) but reduces recall (52 true claims incorrectly rejected).

H Use of LLM

We have only used LLM for language polishing purposes in the paper writing. We do not use LLM for idea formalization, or to an extent that it could be regarded as a contributor.

Category	Count	Description
RULES over-rejection	52	Target=True, RULES=False, baseline correct.
RULES under-rejection	9	Target=False, RULES=True, baseline correct.
RULES no valid answer	1	RULES fails to produce a valid True/False answer, baseline correct.
RULES correct rejection	26	Target=False, baseline=True, RULES=False.
RULES correct acceptance	7	Target=True, baseline=False, RULES=True.
Both correct	1025	Baseline and RULES both answer correctly.
Both wrong	96	Baseline and RULES both answer incorrectly.

Table 3: **CounterMath precision/recall tradeoff for DeepMath-1.5B.** RULES errors are much more likely to be over-rejections than under-rejections, reflecting a stronger tendency to reject claims whose validity it cannot confidently verify.