# LongCat

# AMO-Bench: Large Language Models Still Struggle in High School Math Competitions

**Shengnan An**[*◇]**, Xunliang Cai**[◇]**, Xuezhi Cao**[*◇]**, Xiaoyu Li**[◇]**, Yehao Lin**[◇]**, Junlin Liu**[†♣]**,**
**Xinxuan Lv**[◇]**, Dan Ma**[◇]**, Xuanlin Wang**[†♡]**, Ziwen Wang**[◇]**, Shuang Zhou**[◇]
(Alphabetical order by last name)

[◇]Meituan   [♣]University of Chinese Academy of Sciences   [♡]Harbin Institute of Technology

## ABSTRACT

We present AMO-Bench, an **A**dvanced **M**athematical reasoning benchmark with **O**lympiad level or even higher difficulty, comprising 50 human-crafted problems. Existing benchmarks have widely leveraged high school math competitions for evaluating mathematical reasoning capabilities of large language models (LLMs). However, many existing math competitions are becoming less effective for assessing top-tier LLMs due to performance saturation (e.g., AIME24/25). To address this, AMO-Bench introduces more rigorous challenges by ensuring all 50 problems are (1) cross-validated by experts to meet at least the International Mathematical Olympiad (IMO) difficulty standards, and (2) entirely original problems to prevent potential performance leakages from data memorization. Moreover, each problem in AMO-Bench requires only a final answer rather than a proof, enabling automatic and robust grading for evaluation. Experimental results across 26 LLMs on AMO-Bench show that even the best-performing model achieves only 52.4% accuracy on AMO-Bench, with most LLMs scoring below 40%. Beyond these poor performances, our further analysis reveals a promising scaling trend with increasing test-time compute on AMO-Bench. These results highlight the significant room for improving the mathematical reasoning in current LLMs. We release AMO-Bench to facilitate further research into advancing the reasoning abilities of language models.

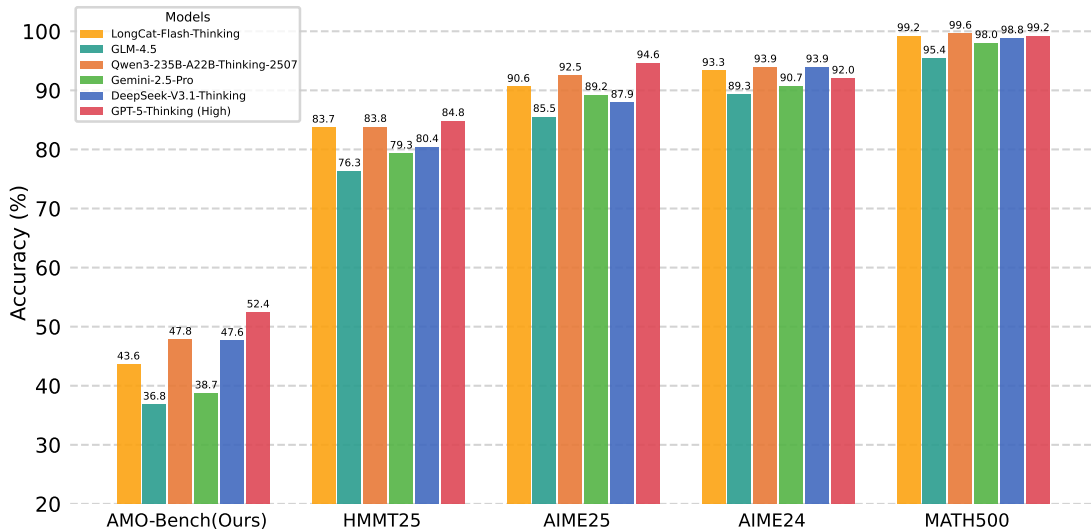**Code, Dataset, and Leaderboard:**  amo-bench.github.io



Figure 1: Performance of top-tier reasoning models on AMO-Bench as well as existing competition-level math benchmarks. Except for the results on AMO-Bench, all other results are sourced from Meituan LongCat Team [2025a].

---

[*] Correspondence to: {anshengnan, caoxuezhi}@meituan.com.
[†] Work done during the internship at Meituan.

# 1 Introduction

Recent advances in large language models (LLMs) have demonstrated significant improvements in reasoning capabilities [OpenAI, 2024, Gemini Team, 2025, OpenAI, 2025, Anthropic, 2025, xAI, 2025, Yang et al., 2025, Guo et al., 2025, DeepSeek-AI, 2025, Meituan LongCat Team, 2025b, GLM-4.5 Team, 2025, ByteDance Seed, 2025, Tencent Hunyuan Team, 2025, Kimi Team, 2025, Meituan LongCat Team, 2025a]. To track this rapid progress, mathematical problem solving has become a critical metric for evaluation, as it inherently demands complex and multi-step reasoning processes to arrive at correct answers. As a result, many current benchmarks utilize problems from high school mathematics competitions (e.g., HMMT and AIME) to assess the reasoning abilities of LLMs [Balunović et al., 2025, He et al., 2024, Gao et al., 2024, Fang et al., 2025]. Recent results indicate that state-of-the-art models are achieving remarkable performances on these benchmarks, with some even surpassing 90% accuracy on competitions like AIME24/25.

However, these impressive results also expose an emerging challenge: many existing mathematics benchmarks are approaching performance saturation and are becoming less effective for assessing further advancements in reasoning capabilities. On the one hand, as LLMs gradually approach or even surpass human-level capabilities in mathematics, some math competitions are becoming less challenging for top-tier models [OpenAI, 2025, DeepSeek-AI, 2025, Yang et al., 2025, Meituan LongCat Team, 2025a]. On the other hand, most current benchmarks are derived from previous competitions, raising concerns about potential data memorization and performance leakage [Sun et al., 2025, Balunović et al., 2025]. While recent efforts have incorporated problems from more difficult and newly held contests such as the International Mathematical Olympiad (IMO), these questions tend to be proof-based and require manual verification by experts [Balunović et al., 2025, Petrov et al., 2025]. This reliance on expert review hinders the implementation of automated scoring processes, leading to inefficiency and inconsistency in large-scale evaluations and result reproductions.

To address these limitations, we present AMO-Bench, an advanced mathematical reasoning benchmark consisting of 50 novel and extremely challenging problems. The core features of AMO-Bench are as follows:

- **Original problems.** To prevent performance leaks from existing resources as much as possible, all problems in AMO-Bench are newly crafted by human experts. Moreover, we conduct a secondary verification to ensure that there are no highly similar problems in existing competitions or online resources.

- **Guaranteed difficulty.** Each problem has undergone rigorous cross-validation by multiple experts to ensure it meets at least the difficulty standards of IMO. We also incorporate an LLM-based difficulty filtering stage to exclude questions that do not present sufficient challenge to current reasoning models.

- **Final-answer based grading.** Each problem in AMO-Bench requires a final answer rather than a full proof, enabling efficient automatic grading. For each problem, we employ a parser-based or LLM-based grading method according to its answer type, balancing the grading cost and generalizability.

- **Human-annotated reasoning paths.** In addition to the final answer, each problem also includes a detailed reasoning path written by human experts. These additional annotations enhance solution transparency and could support further explorations on AMO-Bench, such as prompt engineering and error analysis.

Experimental results across various LLMs demonstrate that contemporary LLMs still struggle with the significant challenges presented by AMO-Bench. Among 26 evaluated models, the state-of-the-art accuracy on AMO-Bench is only 52.4%, achieved by GPT-5-Thinking (High), with most models scoring below 40%. Figure 1 illustrates the performance of several leading models on AMO-Bench as well as the comparison with other mathematical benchmarks. Beyond their limited final performances on AMO-Bench, LLMs consume substantially more output tokens in AMO-Bench compared to existing evaluation datasets. For example, GPT-5-Thinking (High) generates an average of approximately 37K output tokens for AMO-Bench, whereas it produces only about 7K and 6K tokens for AIME25 and AIME24, respectively. This exceptionally high token consumption further underscores the difficulty of AMO-Bench for current LLMs. Despite the poor performances of current LLMs, our analysis also reveals considerable potential for further improvements. Notably, top-tier models achieve pass@32 rates exceeding 70%, suggesting they possess the initial capability to solve these challenging problems even if they do not consistently identify the correct reasoning path at present. Furthermore, we show that the model performances exhibit a near-linear growth trend relative to the logarithm of output length, indicating continued benefits from test-time scaling. These analyses suggest substantial opportunities remain to enhance reasoning capabilities in future generations of language models.

The data and evaluation code of AMO-Bench are publicly available at `amo-bench.github.io`. We hope this novel and challenging benchmark will facilitate further research into advancing the reasoning abilities of language models.
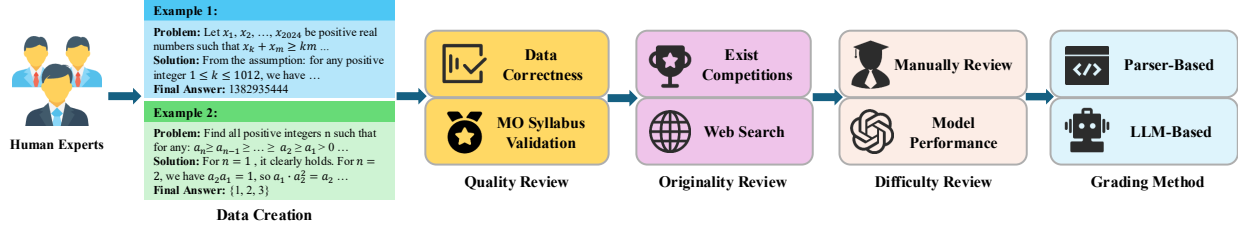
Figure 2: The construction and grading pipeline of AMO-Bench.

## 2 AMO-Bench

In this section, we first introduce the construction process of AMO-Bench (Section 2.1) and present the basic statistics of this dataset (Section 2.2). Then, we elaborate on the grading methodology designed for AMO-Bench (Section 2.3). Figure 2 briefly illustrate the construction and grading pipeline of AMO-Bench.

### 2.1 Construction Pipeline

To ensure the high standards of quality, originality, and difficulty level in our dataset, we have built up a comprehensive multi-stage construction pipeline that covers the entire process from question creation to final inclusion. This pipeline comprises four major stages: data creation, quality review, originality review, and difficulty review.

**Data creation.** All problems are independently designed by mathematics experts from top universities and educational institutions. These experts have extensive backgrounds in high school mathematics competitions, either having won MO-level mathematics competition awards or possessing experience in competition problem design. Beyond the final answer, each problem author must provide a detailed step-by-step solution. These annotated solutions will be utilized in the subsequent quality review stage and will also aid in assessing the overall difficulty of AMO-Bench (see Section 2.2 for details).

**Quality review.** Each candidate problem undergoes blind review by at least three experts to assess its quality. This quality review stage focuses primarily on two aspects:

- Whether the problem statement and solution are semantically unambiguous and logically correct.
- Whether the mathematical knowledge required for the problem is within the scope typically covered in MO-level competitions such as IMO.

**Originality review.** The originality review stage aims to ensure that these newly created problems are not mere rewrites of publicly available materials, but demonstrate genuine originality. To this end, we assess the originality of each problem through the following methods:
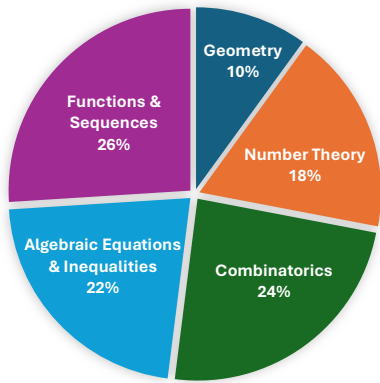
- Compare it against problems in existing datasets (e.g., AIME24/25) with 10-gram matching.
- Conduct web searches to identify any similar online content.

Additionally, during the quality review stage, experts are also required to indicate whether they have encountered highly similar questions in past competitions.
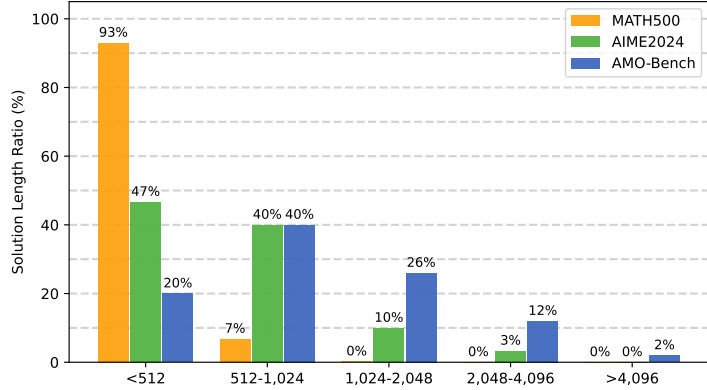
**Difficulty review.** To ensure that AMO-Bench presents a sufficient challenge to state-of-the-art LLMs, we implement a difficulty review stage to filter out problems lacking adequate complexity (even if they may be suitable for some MO-level competitions, e.g., the first 10 questions in AIME). Specifically, each selected problem must satisfy the following two criteria:

- The problem must meet or exceed the IMO difficulty standards, as verified by the human expert.
- We employed multiple advanced reasoning models (such as GPT, DeepSeek, and Gemini series models) for preliminary evaluation, requiring that at least two such models fail to correctly and consistently solve the problem[3].

---

[3]For each model, our preliminary evaluation involves three samples. If all three samples are correct, the model is deemed capable of consistently solving the problem.

(a) Distribution of problem categories.



(b) Comparison of solution lengths.

Figure 3: Basic statistics of AMO-Bench. (a) The distribution of problem categories in AMO-Bench. (b) The distribution of human-annotated solutions in AMO-Bench as well as the comparison with MATH500 and AIME24.

## 2.2 Dataset Statistics

**Problem categories.**    Referring several official competition syllabus, we categorize the 50 problems of AMO-Bench into the following five primary categories: Algebraic Equations & Inequalities (11/50), Functions & Sequences (13/50), Geometry (5/50), Number Theory (9/50), and Combinatorics (12/50). Figure 3a show the overall distribution of problem categories in AMO-Bench.

**Length distribution of human-annotated solutions.**    Since the problems in our AMO-Bench are equipped with manually annotated solutions, we can preliminarily analyze the reasoning complexity of these problems from the view of solution length. We measure solution length in terms of token count[4]. Additionally, we compare the distribution of solution lengths with those from AIME24[5] and MATH500[6]. Figure 3b illustrates the solution length distributions across these benchmarks. It reveals that solutions in AMO-Bench exhibit significantly higher lengths, indicating that problems in this benchmark are inherently more challenging and require more complex reasoning to arrive at the final answer. We conduct a further analysis of the model solution lengths in Section 3.2.

## 2.3 Grading Method

For evaluating answers generated by LLMs, prior work has primarily utilized two approaches: parser-based grading and LLM-based grading. Parser-based grading offers high efficiency and accuracy when the model's response can be successfully parsed; however, its applicability is limited to simple answer formats such as numerical values or sets, making it challenging to assess more complex answers. In contrast, LLM-based grading provides greater flexibility across diverse answer types but may be less efficient and does not consistently guarantee accuracy.

To fully leverage the strengths of both grading methods, AMO-Bench employs different grading approaches based on the specific answer type for each problem. Specifically, problems in AMO-Bench are divided into four main answer types: numerical answers (e.g., Example 1), set answers (e.g., Example 2), variable-expression answers (e.g., Example 3 which requires providing the general formula for an arithmetic sequence), and descriptive answers (e.g., Example 4 which involves comprehensively considering multiple scenarios). The prompt templates for used for grading are contained in Appendix A.

---

**Example 1: Problem with Numerical Answer**

**Question:** Let $x_1, x_2, \cdots, x_{2024}$ be positive real numbers such that $x_k + x_m \geq km$ for any $1 \leq k < m \leq 2024$. Find the minimum value of $x_1 + x_2 + \cdots + x_{2024}$.

**Answer:** $\boxed{1382935444}$

---

[4]We use the tokenizer of DeepSeek-V3.1 model to count tokens in solutions.
[5]https://huggingface.co/datasets/HuggingFaceH4/aime_2024.
[6]https://huggingface.co/datasets/HuggingFaceH4/MATH-500.

---

**Example 2: Problem with Set Answer**

**Question:** Find all positive integers **n** such that for any: $a_n \geq a_{n-1} \geq a_{n-2} \geq \cdots\cdots a_2 \geq a_1 > 0$, satisfying $\sum\limits_{k=1}^{n} a_k = \sum\limits_{k=1}^{n} \frac{1}{a_k}$, the inequality $\prod\limits_{k=1}^{n} a_k^k \geq 1$ holds.

**Answer:** $\boxed{\{1, 2, 3\}}$

---

**Example 3: Problem with Variable-Expression Answer**

**Question:** The sequence $\{a_n\}_{n=1}^{\infty}$ consists of positive terms, with $a_1 = 7$, $a_2 = 2$, and satisfies the recurrence relation

$$8a_{n+2}^4 = 3 + 4a_{n+1} + a_n \quad (n \in \mathbb{N}^*).$$

Find the general term formula for this sequence.

**Answer:**

$$\boxed{\frac{(2 + \sqrt{3})^{2^{2-n}} + (2 - \sqrt{3})^{2^{2-n}}}{2}}$$

---

**Example 4: Problem with Descriptive Answer**

**Question:** Let $n$ be an integer with $n > 2$. Real numbers $a_1, a_2, \ldots, a_n$ satisfy

$$\sum_{k=1}^{n} a_k = 2n, \qquad \sum_{k=1}^{n} k\,|a_k| = 4n.$$

Find the minimum value of $a_1^2 + a_2^2 + \cdots + a_n^2$.

**Answer:** For $n = 3$, the minimum of $a_1^2 + a_2^2 + a_3^2$ is 12.

For $n \geq 4$, the minimum of $a_1^2 + a_2^2 + \cdots + a_n^2$ is $\frac{6n^2}{5}$.

---

For problems requiring numerical, set, or variable-expression answers (39 out of 50), we employ the parser-based grading. The evaluated LLMs are instructed to format their final responses as \boxed{<answer>}. We then utilize the tools provided by math-verify[7] to parse these answers and verify the equivalence with the ground truth. Moreover, if the model answer containing decimal values, we require an accuracy of at least four decimal places. For variable-expression answers, we assign multiple sets of values to the variables in the expression, then verify whether the values of the generated expression match that of the ground-truth expression. We also manually review the parsing results during the preliminary evaluation and adjust the post-processing algorithms.

For problems requiring descriptive answers (11 out of 50), we use LLM-based grading with o4-mini (Low) serving as the grading model. To ensure robust assessment, majority voting is performed across five independent grading samples for each response. Additionally, during preliminary evaluation, we manually verify the correctness of LLM-based grades for all descriptive answers and revise answer descriptions where needed to enhance grading accuracy.

**Grading accuracy.** Prior to conducting the large-scale evaluation, we performed a manual quality check to ensure the reliability of the designed grading method. This assessment included 1,000 responses generated by 10 different LLMs. The results indicate that the grading accuracy reached 99.2%, providing strong validation for the effectiveness of the grading method on AMO-Bench.

## 3 Experiments

In this section, we present the experimental results on AMO-Bench. We first describe the experimental setup (Section 3.1), followed by a discussion of the main results and analysis (Section 3.2).
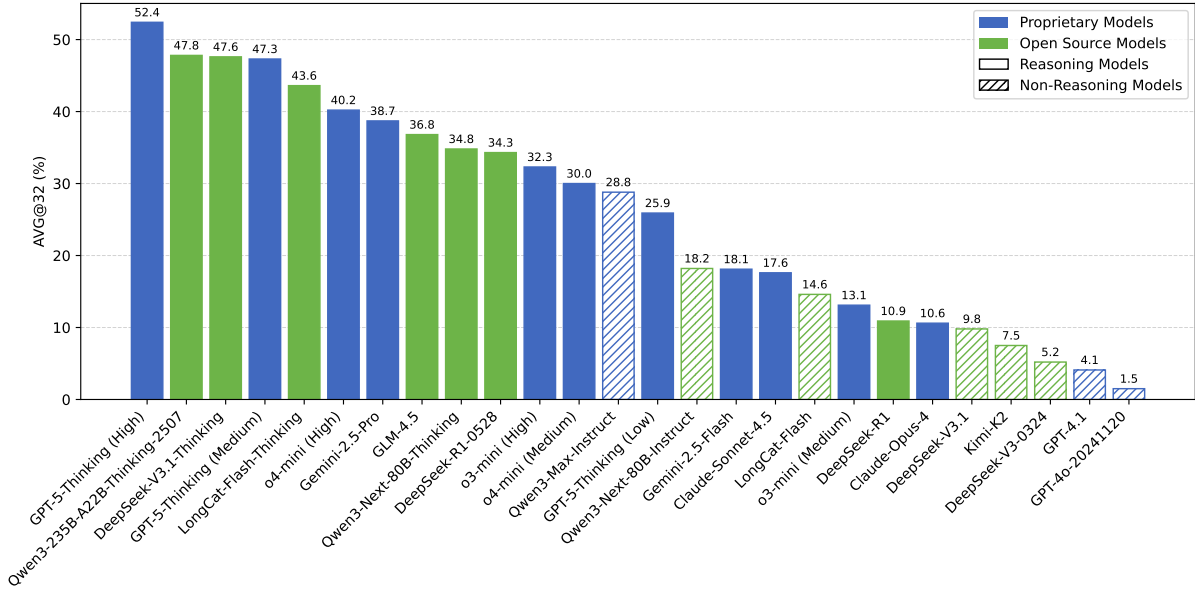
---

[7]https://github.com/huggingface/Math-Verify.

Figure 4: The AVG@32 performance of various LLMs on AMO-Bench.

## 3.1 Experimental Setup

**Models.** To conduct a comprehensive and representative evaluation on AMO-Bench, we select a diverse set of leading LLMs, encompassing both open-source models and proprietary models. Specifically, the evaluation includes top-tier models provided by OpenAI [OpenAI, 2025], Gemini [Gemini Team, 2025], Anthropic [Anthropic, 2025], DeepSeek [Guo et al., 2025], Qwen [Yang et al., 2025], GLM [GLM-4.5 Team, 2025], Moonshot [Kimi Team, 2025], and LongCat [Meituan LongCat Team, 2025a]. In addition to evaluating reasoning models that have been specifically enhanced for long-term thinking tasks, we also incorporated several powerful non-reasoning models to demonstrate their potential in tackling complex reasoning challenges.

**Sampling settings.** We set the `temperature` of sampling to 1.0 for reasoning models and 0.7 for non-reasoning models. For all evaluated models, we use `top-k=50` and `top-p=0.95` during sampling. We configure the maximum context/output length to the highest allowable limit for each model during inference. This avoids underestimating the reasoning capabilities of the model due to restrictions on the token budget. To ensure the stability of the final evaluation results, we sampled the results from each model 32 times and reported the average performance of these 32 results as the final metric (denoted as AVG@32). Appendix B illustrates the fluctuation of the average result across different sampling times. It demonstrates that when sampling 32 times, the average model performance exhibits a relatively small fluctuation and rarely appears to reverse the model ranking order.

## 3.2 Results and Analysis

**Main results.** Figure 4 presents the AVG@32 performance of various leading LLMs, categorized by proprietary/open-source status and reasoning/non-reasoning properties[8]. Overall, all these models still struggle with the significant challenges presented by AMO-Bench. Even the highest performing model GPT-5-Thinking (High) reaches just 52.4%, while most others score below 40%. This indicates substantial room for improvement in complex reasoning abilities across all current language models. Moreover, both proprietary and open-source reasoning models occupy top ranks in the leaderboard, indicating that recent open-source advancements are closing the gap with leading commercial models. The best-performing open-source model is only about 5% lower than the top proprietary result. Besides reasoning models, some non-reasoning models demonstrate a performance exceeding expectations, such as Qwen3-Max-Instruct

---

[8]To facilitate easier reproduction and utilization of AMO-Bench, you can take a fast try on the AMO-Bench-P subset, which includes only the 39 parser-based grading problems from AMO-Bench. Appendix C presents the AVG@32 performance of LLMs on AMO-Bench-P.
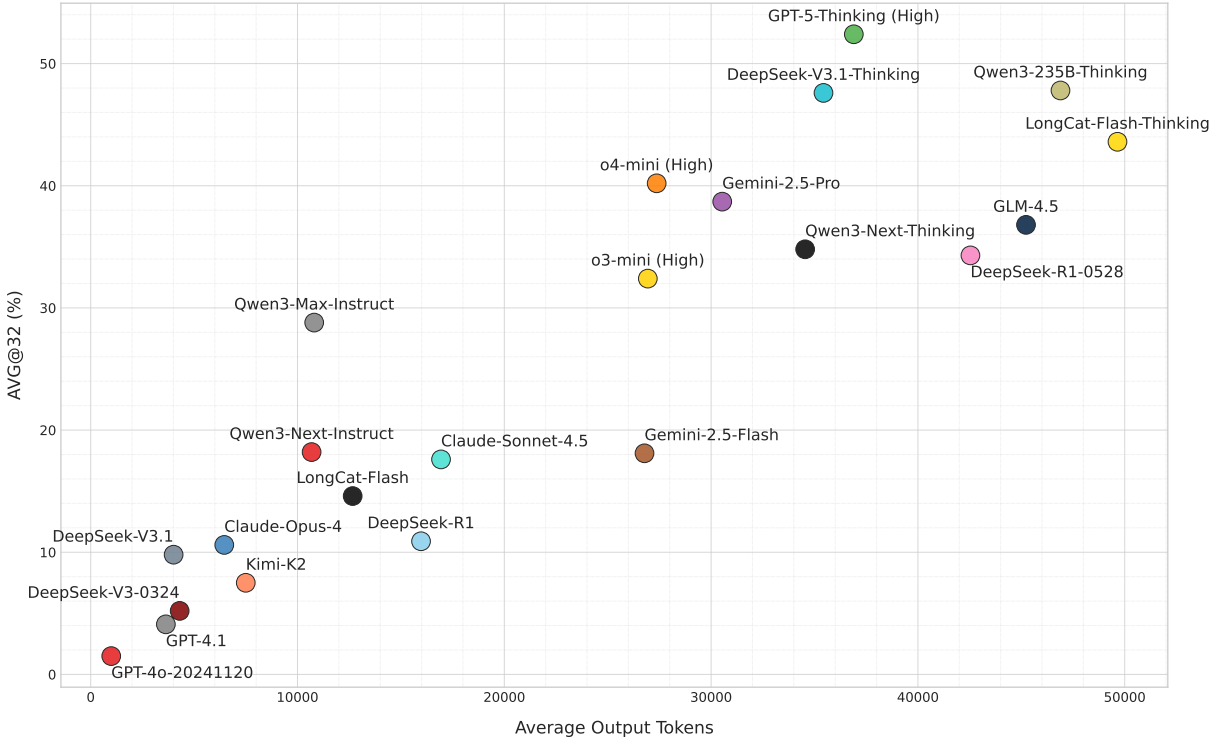
Figure 5: The AVG@32 performance of LLMs vs. the average model output length.

and LongCat-Flash. These non-reasoning models even outperforms several reasoning models such as o3-mini (Medium), indicating their significant potential in tackling complex reasoning tasks.

**Comparison of reasoning efficiency.**  Figure 5 shows the average output length and the AVG@32 performance of each model. Overall, it demonstrates a clear trend that higher-performing models tend to require more output tokens. The first-tier models that reach higher than 40% AVG@32 scores utilize more than 35K completion tokens. Even among non-reasoning models, those with superior performance are distinguished by their ability to process more tokens, sometimes reaching levels comparable to reasoning models. Additionally, when examining models within the same series, there are notable improvements in reasoning efficiency over time. For example, o4-mini (High) outperforms o3-mini (High) at similar or slightly increased token counts. Likewise, DeepSeek-V3.1-Thinking shows significant gains compared to DeepSeek-R1-0528 with even significantly less output tokens.

Beyond the main results outlined above, we also provide further analysis and insights based on the AMO-Bench experimental findings.

**The model output length could indicate the reasoning challenge of the benchmark.**  Section 2.2 provides a pre-analysis of benchmark difficulty based on annotated solution lengths. Here, we offer a post-hoc analysis of benchmark difficulty based on the relationship between model performance and model output length. Figure 6 clearly demonstrates that the average output length of each model increases as the reasoning benchmark becomes more challenging. Specifically, across six models, benchmarks with higher accuracy scores (such as MAH500 and AIME24) correspond to shorter average outputs, while those with lower scores (like AMO-Bench) require significantly longer responses. This suggests that harder benchmarks demand more elaborate reasoning steps or explanations from the models, resulting in increased token usage. These results demonstrate that the model output length could be an indicator of reasoning challenge in the benchmark.

**Performance on AMO-Bench still benefits from test-time scaling.**  The reasoning efficiency results discussed above indicate a correlation between model performance and output length. Here, we conduct a more rigorous analysis by directly controlling the reasoning effort for the same model. As shown in the Figure 7, all three models (GPT-5, o4-mini, and o3-mini) exhibit a near-linear growth trend in AVG@32 as the logarithm of average output length increases. Such a trend is highly aligned with earlier experimental observations from existing benchmarks such as MATH500
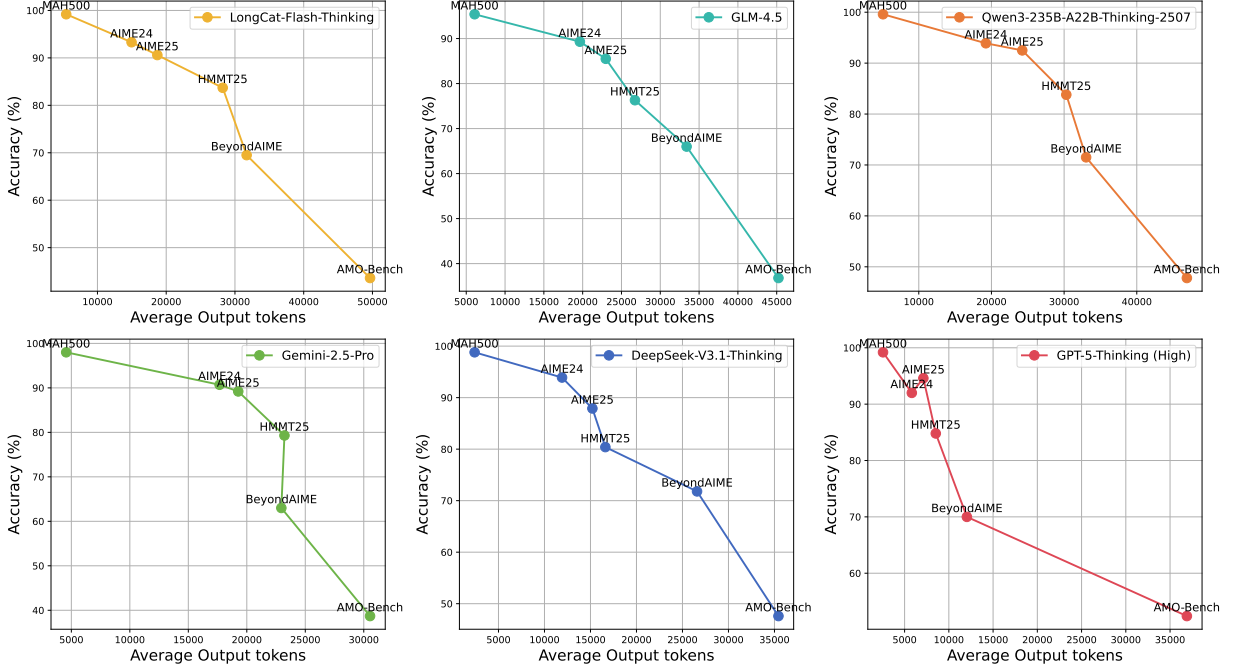
Figure 6: The relationship between accuracy and average output length on different math benchmarks.
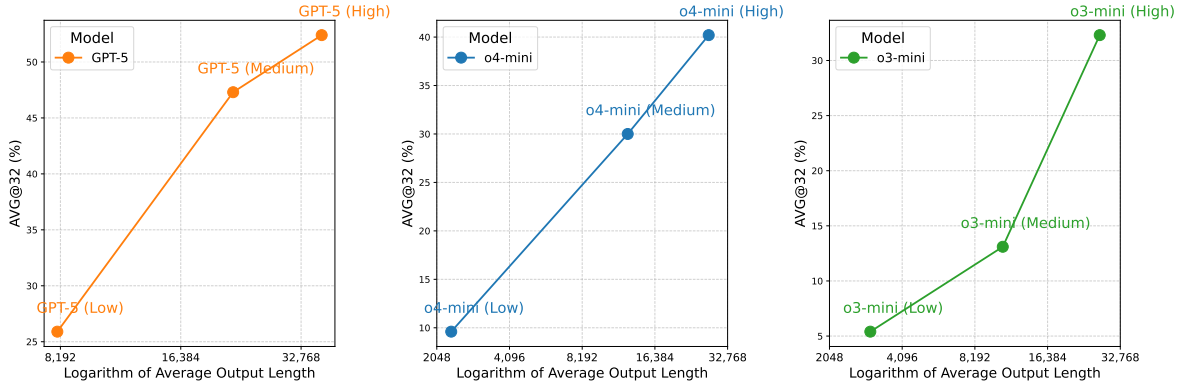


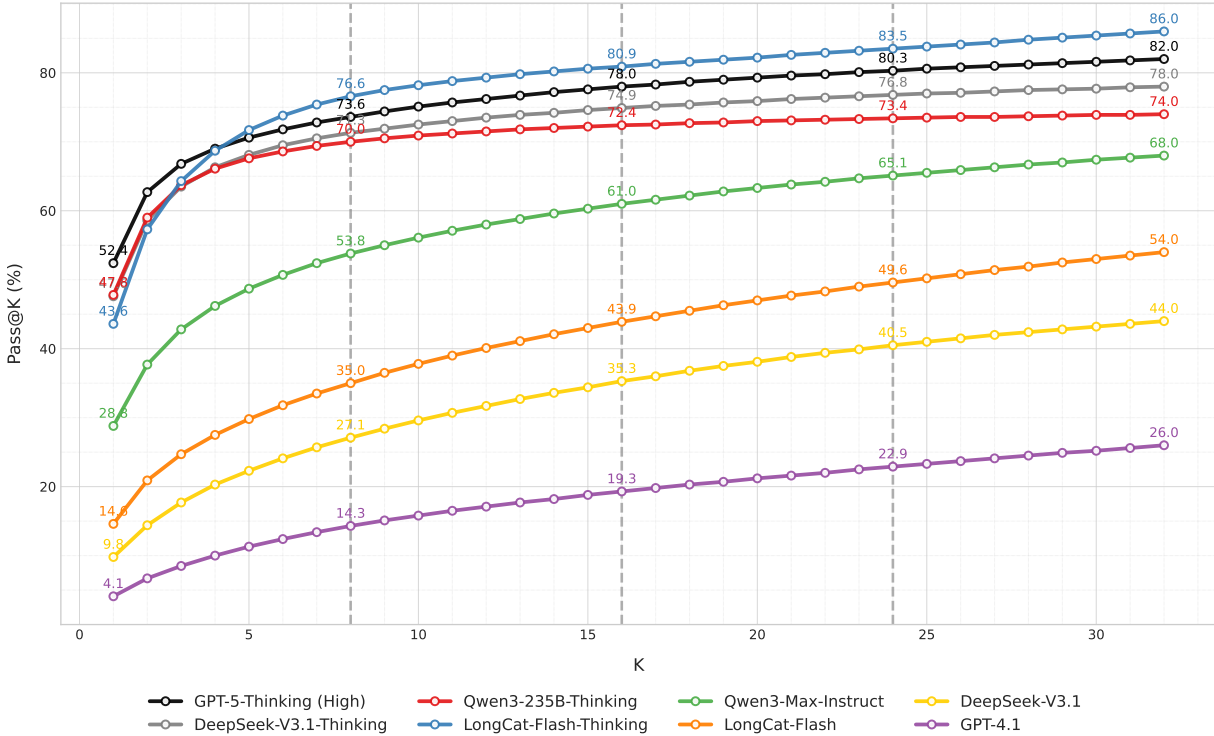Figure 7: The model performance and output length under different reasoning effort settings.

and AIME24 [Muennighoff et al., 2025]. This indicates that further increasing the inference budget will further drive improvements on AMO-Bench.

**Top-tier models demonstrate promising potential for improvement on AMO-Bench.** Existing work reveals that the pass@$k$ performance of the model can reflect its inherent potential to achieve further improvement through reinforcement learning. Inspired by this, we illustrate the pass@$k$ of evaluated models to indicate their inner potential. As shown in Figure 8, the pass@$k$ metric exhibits rapid growth as $k$ increases from 1 to 8, followed by a sustained but gradual improvement as $k$ continues to rise. Notably, the top-tier reasoning models achieve over 70% performance on the pass@32 metric. These results highlight the significant room for improvement in the reasoning capabilities of LLMs.

## 4  Related Work

Evaluating LLMs on mathematical problem solving has been a critical aspect for assessing advancements in reasoning capabilities. Early datasets such as GSM8K [Cobbe et al., 2021] and MATH [Hendrycks et al., 2021] provided

Figure 8: The the pass@$k$ trend of various LLMs with increasing $k$.

initial explorations to evaluate these abilities. However, model performance on these benchmarks has quickly reached saturation. To further advance the study of mathematical proficiency in LLMs, recent work has shifted toward more challenging benchmarks.

In terms of increasing difficulty, two primary lines of work have emerged. One line focuses on Mathematical Olympiad (MO)-level problems, which rely on a specific range of math knowledge and require complex and intuitive reasoning skills. For instance, Omni-MATH [Gao et al., 2024] introduces a multi-subject evaluation suite designed to rigorously test mathematical reasoning and generalization; OlympiadBench [He et al., 2024] focuses on evaluating the bilingual and multi-modal reasoning abilities with Olympid-level challenges; OlymMATH [Sun et al., 2025] collects MO-level problems from printed publications and evaluates mathematical reasoning by offering problems of two difficulty levels; MathOdyssey [Fang et al., 2025] broadens the scope to include more complex tasks, with a particular focus on long-range and compositional reasoning; BeyondAIME [ByteDance-Seed, 2025] collects problems similar in style to AIME with increased difficulty and expanded data scale; MathArena [Balunović et al., 2025] rapidly tracks model performance in newly held MO-level competitions and explores evaluation paradigms for proof-based competitions such as the IMO and USAMO. Our proposed AMO-Bench also falls within this category and it stands as one of the most challenging benchmarks at the time of writing.

The other line of work focuses on problems derived from graduate-level examinations or advanced mathematical research. For instance, RealMath [Zhang et al., 2025] provides a comprehensive evaluation of LLMs in real-world mathematical tasks, assessing their reasoning capabilities across a diverse range of research-level content; Frontier-Math [Glazer et al., 2024] covers computationally intensive problems and abstract questions across most branches of mathematics, highlighting the significant gap between LLMs and the prowess of the mathematical community; HARDMath2 [Roggeveen et al., 2025] focuses on approximation-based mathematical problems, particularly those commonly encountered in applied sciences and engineering; HLE [Phan et al., 2025] constructs a final closed-ended academic benchmark spanning multiple subjects, evaluating reasoning capabilities on human frontier knowledge. Beside requiring the reasoning abilities, these datasets also challenge models by demanding extensive and deep mathematical knowledge.

# 5    Conclusion

We introduce AMO-Bench, an advanced mathematical reasoning benchmark featuring problems at the level of mathematical Olympiads or higher. The benchmark consists of 50 human-crafted questions designed to rigorously assess advanced mathematical reasoning. Compared with existing benchmarks, AMO-Bench offers more challenging assessments by ensuring that all 50 problems are entirely original and meet or exceed IMO difficulty standards. Each problem in AMO-Bench requires only a final answer rather than a full proof, enabling automatic and robust grading for evaluation purposes. Experimental results across various LLMs demonstrate that contemporary LLMs still struggle with the significant challenges presented by AMO-Bench. Despite these low performances, our further analysis underscore substantial opportunities for advancing mathematical reasoning capabilities in current LLMs.

# Acknowledgments

# References

Meituan LongCat Team. Longcat-flash-thinking technical report, 2025a. URL https://arxiv.org/abs/2509.18883.

OpenAI. Openai o1 system card, 2024. URL https://arxiv.org/abs/2412.16720.

Gemini Team. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. URL https://arxiv.org/abs/2507.06261.

OpenAI. Gpt-5 system card, 2025. URL https://cdn.openai.com/gpt-5-system-card.pdf.

Anthropic. System card: Claude opus 4 and claude sonnet 4, 2025. URL https://www-cdn.anthropic.com/4263b940cabb546aa0e3283f35b686f4f3b2ff47.pdf.

xAI. Grok 4 model card, 2025. URL https://data.x.ai/2025-08-20-grok-4-model-card.pdf.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. Nature, 645(8081): 633–638, 2025.

DeepSeek-AI. Deepseek-v3 technical report, 2025. URL https://arxiv.org/abs/2412.19437.

Meituan LongCat Team. Longcat-flash technical report, 2025b. URL https://arxiv.org/abs/2509.01322.

GLM-4.5 Team. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models, 2025. URL https://arxiv.org/abs/2508.06471.

ByteDance Seed. Seed1.5-thinking: Advancing superb reasoning models with reinforcement learning, 2025. URL https://arxiv.org/abs/2504.13914.

Tencent Hunyuan Team. Hunyuan-turbos: Advancing large language models through mamba-transformer synergy and adaptive chain-of-thought, 2025. URL https://arxiv.org/abs/2505.15431.

Kimi Team. Kimi k2: Open agentic intelligence, 2025. URL https://arxiv.org/abs/2507.20534.

Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. Matharena: Evaluating llms on uncontaminated math competitions. arXiv preprint arXiv:2505.23281, 2025.

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3828–3850, 2024.

Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, et al. Omni-math: A universal olympiad level mathematic benchmark for large language models. In The Thirteenth International Conference on Learning Representations, 2024.

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. Scientific Data, 12(1):1392, 2025.

Haoxiang Sun, Yingqian Min, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, Zheng Liu, Zhongyuan Wang, and Ji-Rong Wen. Challenging the boundaries of reasoning: An olympiad-level math benchmark for large language models. arXiv preprint arXiv:2503.21380, 2025.

Ivo Petrov, Jasper Dekoninck, Lyuben Baltadzhiev, Maria Drencheva, Kristian Minchev, Mislav Balunović, Nikola Jovanović, and Martin Vechev. Proof or bluff? evaluating llms on 2025 usa math olympiad, 2025. URL https://arxiv.org/abs/2503.21934.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candes, and Tatsunori Hashimoto. s1: Simple test-time scaling. In Workshop on Reasoning and Planning for Large Language Models, 2025.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.

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

ByteDance-Seed. Beyondaime: Advancing math reasoning evaluation beyond high school olympiads, 2025. URL https://huggingface.co/datasets/ByteDance-Seed/BeyondAIME.

Jie Zhang, Cezara Petrui, Kristina Nikolić, and Florian Tramèr. Realmath: A continuous benchmark for evaluating language models on research-level mathematics. arXiv preprint arXiv:2505.12575, 2025.

Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caroline Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, et al. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai. arXiv preprint arXiv:2411.04872, 2024.

James V Roggeveen, Erik Y Wang, Will Flintoft, Peter Donets, Lucy S Nathwani, Nickholas Gutierrez, David Ettel, Anton Marius Graf, Siddharth Dandavate, Arjun Nageswaran, et al. Hardmath2: A benchmark for applied mathematics built by students as part of a graduate class. arXiv preprint arXiv:2505.11774, 2025.

Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity's last exam. arXiv preprint arXiv:2501.14249, 2025.

# A    Prompt Templates

**Query prompt template.**    In order to guide LLMs in generating answers in a parser-readable format, we use the following prompt template guide the model generation. There are mainly three requirements in the instruction: the answer prefix (i.e., ### The final answer is:), the LaTeX box environment (i.e., \boxed{}), and the precision requirement.

---

**Example 5: Query Prompt Template**

...
After solving the above problem, please output your final answer in the following format:
### The final answer is: $\boxed{<your answer>}$
Example:
### The final answer is: $\boxed{123}$
The final answer should be given as precisely as possible (using LaTeX symbols such as \sqrt, \frac, \pi, etc.). If the final answer involves a decimal approximation, it must be accurate to at least four decimal places.

---

**Grading prompt template.**    We employ the LLM-based grading using o4-mini (Low) as the grading model, and use the following grading prompt to verify the equivalence between the LLM output and the reference answer.

---

**Example 6: Grading Prompt Template**

For the following math problem, we have the reference answer and the student's answer.
Determine whether the student's answer is equivalent to the reference answer.
If equivalent, output "Correct".
If not equivalent, output "Incorrect".

### Problem
...

### Reference Answer
...

### Student Answer
...

Now, please provide your judgment.
Please strictly follow the format below to summarize your conclusion at the end of your judgment:
### Conclusion: Correct/Incorrect
If the answer involves a decimal approximation, it must be accurate to at least four decimal places.

---

# B    Analysis of AVG@$k$

Figure 9 illustrates the fluctuation of the average performance across different sampling times. It shows that as the sampling time grows, the models' performance become more stable. When sampling 32 times, it rarely appears the reverse-order phenomenon.

# C    Performance on AMO-Bench-P Subset

To facilitate easier reproduction and use of AMO-Bench, you can utilize the AMO-Bench-P subset, which includes only the 39 parser-based grading problems from AMO-Bench. Table 1 presents the AVG@32 performance of LLMs on AMO-Bench-P. In general, performance on AMO-Bench-P tends to be slightly higher than on the full AMO-Bench, as problems requiring complex descriptive answers are inherently more challenging than those with simple-format answers.
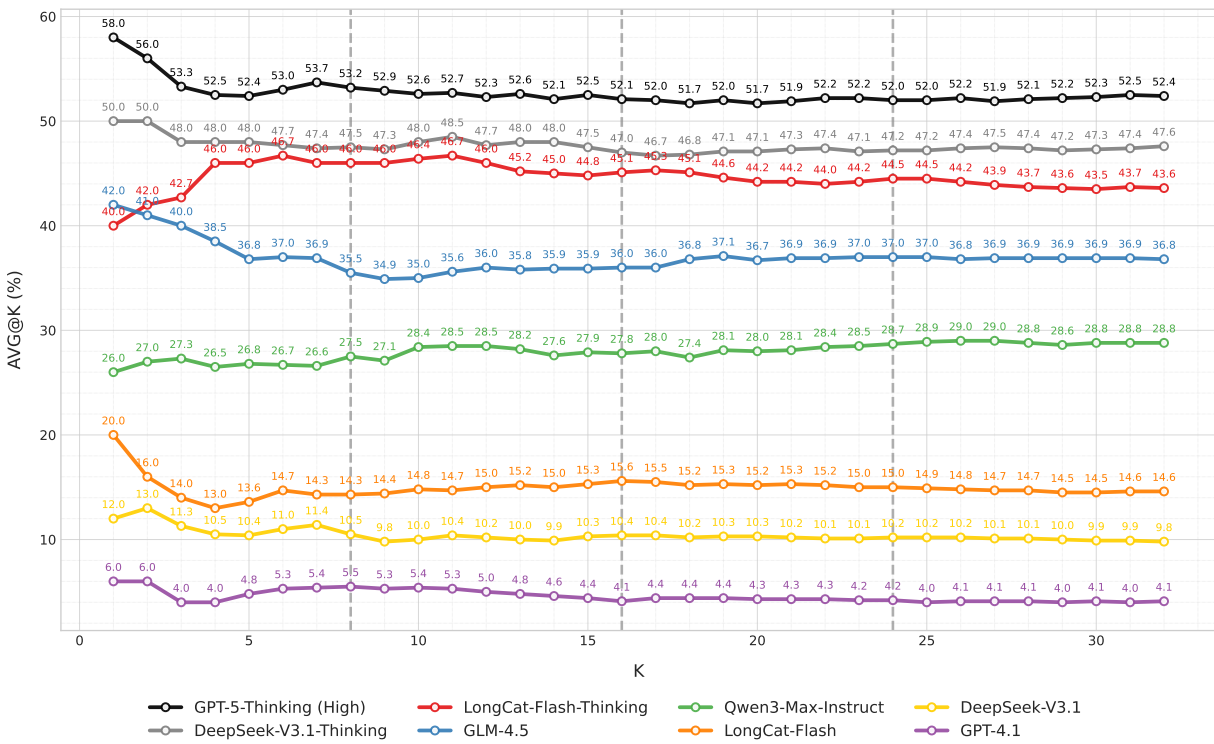
Figure 9: The AVG@$k$ trend of various LLMs with increasing $k$.

Table 1: The AVG@32 performance of LLMs on the AMO-Bench and AMO-Bench-P, the latter of which contains only 39 parser-based grading problems.

| Model | AMO-Bench | AMO-Bench-P |
|---|---|---|
| GPT-5-Thinking (High) | 52.4 | 54.8 |
| Qwen3-235B-A22B-Thinking-2507 | 47.8 | 56.2 |
| DeepSeek-V3.1-Thinking | 47.6 | 53.0 |
| LongCat-Flash-Thinking | 43.6 | 45.3 |
| o4-mini (High) | 40.2 | 43.8 |
| Gemini-2.5-Pro | 38.7 | 41.7 |
| GLM-4.5 | 36.8 | 41.0 |
| Qwen3-Next-80B-Thinking | 34.8 | 37.4 |
| DeepSeek-R1-0528 | 34.3 | 37.1 |
| o3-mini (High) | 32.3 | 34.0 |
| Qwen3-Max-Instruct | 28.8 | 30.9 |
| Qwen3-Next-80B-Instruct | 18.2 | 17.8 |
| Gemini-2.5-Flash | 18.1 | 18.0 |
| Claude-Sonnet-4.5 | 17.6 | 18.1 |
| LongCat-Flash | 14.6 | 14.9 |
| DeepSeek-R1 | 10.9 | 11.7 |
| Claude-Opus-4 | 10.6 | 11.4 |
| DeepSeek-V3.1 | 9.8 | 9.6 |
| Kimi-K2 | 7.5 | 8.4 |
| DeepSeek-V3-0324 | 5.2 | 5.4 |
| GPT-4.1 | 4.1 | 4.8 |
| GPT-4o-20241120 | 1.5 | 1.9 |