

LLMEval-Fair: A Large-Scale Longitudinal Study on Robust and Fair Evaluation of Large Language Models

Anonymous ACL submission

Abstract

Existing evaluation of Large Language Models (LLMs) on static benchmarks is vulnerable to data contamination and leaderboard overfitting, critical issues that obscure true model capabilities. To address this, we introduce LLMEval-Fair, a framework for dynamic evaluation of LLMs. LLMEval-Fair is built on a proprietary bank of 220k graduate-level questions, from which it dynamically samples unseen test sets for each evaluation run. Its automated pipeline ensures integrity via contamination-resistant data curation, a novel anti-cheating architecture, and a calibrated LLM-as-a-judge process achieving 90% agreement with human experts, complemented by a relative ranking system for fair comparison. An 30-month longitudinal study of nearly 60 leading models reveals a performance ceiling on knowledge memorization and exposes data contamination vulnerabilities undetectable by static benchmarks. The framework demonstrates exceptional robustness in ranking stability and consistency, providing strong empirical validation for the dynamic evaluation paradigm. LLMEval-Fair offers a robust and credible methodology for assessing the true capabilities of LLMs beyond leaderboard scores, promoting the development of more trustworthy evaluation standards.

1 Introduction

The rapid advancement of Large Language Models (LLMs) has led to a proliferation of benchmarks designed to assess their capabilities (Chang et al., 2023; Laskar et al., 2024; Liu et al., 2025). However, these benchmarks predominantly rely on a static evaluation paradigm where models are tested on fixed, publicly available datasets (Chen et al., 2025). This approach is fundamentally vulnerable to data contamination and test set overfitting, contributing to a growing “evaluation crisis” where benchmark scores may no longer reliably reflect a model’s generalizable abilities (Banerjee et al.,

2024; Deng et al., 2024a; Dekoninck et al., 2024). For example, GPT-4 achieved exact match rates of 52% and 57% when guessing the masked options in MMLU (Hendrycks et al., 2021a) test sets, far exceeding random chance (Deng et al., 2024b). Similarly, Qwen-1.8B has been shown to exactly replicate complete n-grams from both training and test splits of GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021b), including 25 exact 5-gram matches in the MATH test set, indicating potential undetected data leakage and memorization (Xu et al., 2024b). Furthermore, static benchmark designs, dataset contamination, and biased evaluation protocols could create misleading perceptions of LLM capabilities, undermining the reliability of current performance assessments (Banerjee et al., 2024).

The crisis compels us to shift our focus from what capabilities to evaluate, such as knowledge and reasoning, to a more foundational question: how to evaluate in a manner that is robust, fair, and resistant to strategic manipulation. Based on our analysis of the current evaluation crisis, we identify three fundamental challenges for constructing a trustworthy evaluation framework. These challenges correspond to the core stages of any assessment: the data, the protocol, and the ranking system.

Challenge 1: How can we ensure the integrity of the evaluation data? The cornerstone of any benchmark is its test set. If the data are public or predictable, they become susceptible to leakage into model training corpora, leading to inflated scores that reflect memorization rather than true capability. Therefore, building a benchmark that is inherently resilient to data contamination is the primary challenge. To address this, we construct a private question bank of over 220k graduate-level questions, augmented with structural variations to mitigate memorization.

Challenge 2: How can we design an unpre-

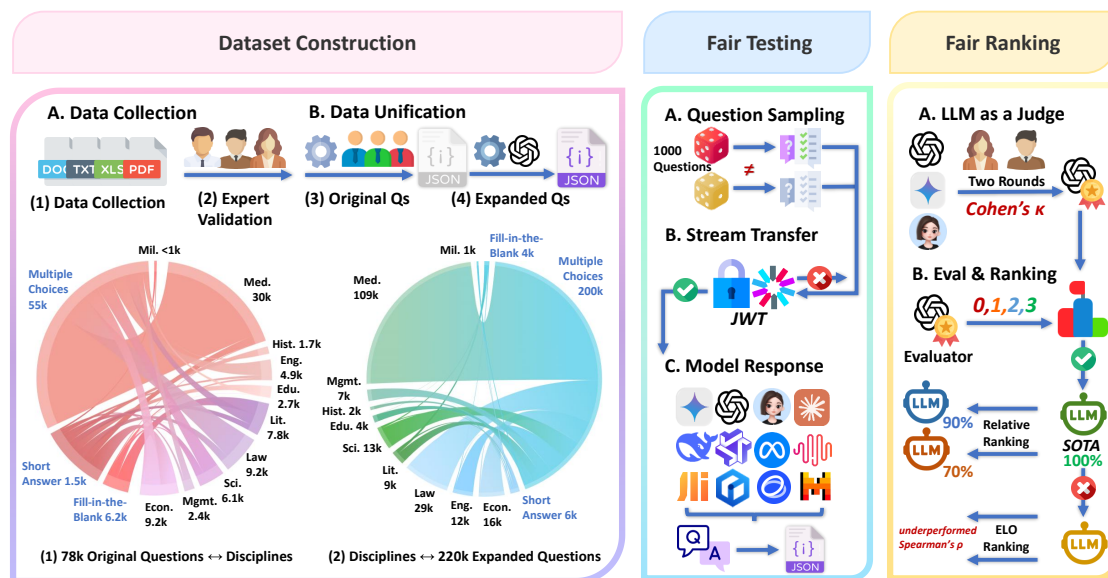


Figure 1: The LLMEval-Fair framework comprises three core stages. First, in Data Construction, diverse exam data are collected, filtered by experts, standardized into original questions in JSON format, and expanded to enhance formality and diversity. Second, Fair Testing involves sampling 1,000 unique questions for models’ evaluation, transmitting data via a secure JWT stream, and collecting model responses. Third, Fair Ranking begins by selecting evaluators based on human-machine agreement, measured by Cohen’s κ coefficients. Subsequently, the process entails generating relative rankings, validating robustness through ablation studies, and analyzing model errors. Discipline abbreviations: Eng. (Engineering), Econ. (Economics), Edu. (Education), Lit. (Literature), Mgmt. (Management), Sci. (Science), Hist. (History), Med. (Medicine), Mil. (Military).

084 **dictable assessment protocol?** Even with pri- 107
 085 vate data, a static protocol with fixed test sets can 108
 086 be reverse-engineered or exploited over time. A 109
 087 truly robust evaluation requires a dynamic and un- 110
 088 predictable process that prevents strategic games. 111
 089 To achieve this, we implement a dynamic evalua- 112
 090 tion protocol where models are served unseen, ran- 113
 091 domly sampled questions in each session through 114
 092 a secure, two-layer anti-cheating architecture, ensur- 115
 093 ing that each evaluation is unique and unpre- 116
 094 dictable. 117

095 **Challenge 3: How can we establish a fair** 118
 096 **and stable ranking system under dynamic con-** 119
 097 **ditions?** When the evaluation content is not fixed, 120
 098 traditional absolute scoring becomes unreliable for 121
 099 cross-model comparison. A fair ranking system 122
 100 must remain stable and consistent, even when mod- 123
 101 els are tested on different, albeit equivalent, sets of 124
 102 questions. To solve this, we develop a novel rela- 125
 103 tive ranking system powered by a highly-calibrated 126
 104 LLM-as-a-Judge framework. This system ranks 127
 105 models by comparing their performance within 128
 106 each evaluation session rather than relying on abso-

lute scores, ensuring fair rankings with negligible 107
 variance across different data samples. 108

109 Collectively, these solutions constitute 110
 LLMEval-Fair, a comprehensive framework for 111
 dynamic LLM evaluation. We leverage this frame- 112
 work to conduct an extensive longitudinal study on 113
 our official platform, spanning from the first half 114
 of 2023 to the second half of 2025. Throughout 115
 this period, we have continuously evaluated 116
 nearly 50 leading proprietary and open-source 117
 models, typically assessing them within one week 118
 of their public release. To ensure the fairness 119
 and timeliness of our evaluation, we constantly 120
 update our private question bank. Each model 121
 undergoes at least three independent, randomly 122
 sampled evaluation runs to ensure the stability of 123
 our findings. To date, this rigorous, ongoing effort 124
 has accumulated over 180k evaluation data points. 125

126 Leveraging this large-scale evaluation data, we 127
 128 systematically investigate three research questions 129
 corresponding to our core challenges and find that 130
 (a) all models converge to a performance ceiling 131
 around 90% persistent gaps in specialized domains 132

like literature and medicine, (b) dynamic rankings diverge from static benchmarks, and static benchmarks suffer from severe data contamination, (c) our framework demonstrates exceptional ranking stability with negligible variance under multi-round resampling and varying sample sizes. More insightful findings are presented in Section 3.

Overall, our contributions are threefold:

1. We construct LLMEval-Fair, a large-scale anti-cheating evaluation platform featuring a proprietary 220k question bank, secure dynamic sampling protocols, and robust anti-manipulation mechanisms for trustworthy LLM assessment.
2. We conduct an extensive 30-month longitudinal evaluation campaign across nearly 60 leading models, accumulating over 180k evaluation data points through continuous anti-cheating assessments and comparative analysis with static benchmarks.
3. We conduct extensive empirical analysis across three research questions corresponding to our core challenges, revealing eight key findings about model performance ceilings, ranking stability, and contamination vulnerabilities in current evaluation practices.

2 Design

In this section, we detail the design and implementation of LLMEval-Fair, which addresses the three fundamental challenges through a three-stage framework. *Dataset Construction* tackles data integrity by building a contamination-resistant private question bank. *Evaluation Process* addresses unpredictable assessment through dynamic sampling and anti-cheating mechanisms. *Ranking System* implements a calibrated ranking system for fair model comparison.

2.1 Dataset Construction

To build a contamination-resistant and high-quality question bank, we sourced postgraduate and undergraduate exam questions from Chinese universities, covering 13 primary and over 50 secondary academic disciplines. Figure 1 offers a comprehensive overview of the design, highlighting three key stages: Dataset Construction, Fair Testing and Fair Ranking.

The construction follows a rigorous pipeline. First, we collect original exam questions from diverse formats and invite over 30 expert annotators for quality screening, which yields 78,009 high-quality original questions after eliminating those with factual errors or irrelevant answers. Second, we employ an LLM-driven augmentation process to expand coverage and diversity. For instance, each Multiple-Choice question with n options is converted into n Fill-in-the-Blank variants, while Material Analysis questions are decomposed into multiple true/false questions based on key information. Finally, all augmented questions undergo format verification and metadata enrichment to ensure quality and traceability.

As of early 2025, this process has resulted in the **LLMEval-Fair dataset**, which comprises over 220k questions across six main categories. To maintain evaluation freshness and prevent contamination, we continuously expand the question bank through the same manual collection and automated augmentation pipeline. Detailed statistical analysis of the dataset, including disciplinary breakdown and content diversity, is provided in Appendix A.

2.2 Evaluation Process

To ensure a reliable and fair evaluation, we design a dynamic process centered on a multi-layered anti-cheating architecture. This approach guarantees that each evaluation is unique, robust against manipulation, and accurately reflects a model’s capabilities. The process is built upon two core strategies: dynamic question sampling and a secure delivery architecture.

2.2.1 Dynamic Question Sampling

To ensure unpredictability, each model evaluation is based on a unique set of 1,000 questions randomly sampled from our private question bank. A dynamic algorithm generates a unique, non-repeatable question sequence for each evaluation session, making it impossible to predict upcoming questions based on historical patterns. Furthermore, models are required to answer questions in the pre-allocated order, preventing any “cherry-picking” strategies. This dynamic sampling ensures that every evaluation is a distinct and isolated event, reflecting the model’s true generalization ability rather than its performance on a fixed test set.

2.2.2 Secure Anti-Cheating Architecture

The evaluation process is protected by a secure, two-layer anti-cheating architecture.

- **The Outer Layer (Access Control):** This layer manages authentication and authorization. We use JSON Web Tokens (JWT)(Jones et al., 2015) to secure every API request, ensuring that only authenticated models can participate in an evaluation session. A strict Role-Based Access Control (RBAC) system prevents any cross-session or cross-user data access, isolating each evaluation.
- **The Inner Layer (Process Control):** This layer enforces the evaluation rules. A multi-level quota system tracks the number of questions allocated, pending, and completed, effectively preventing models from attempting to acquire more questions than permitted or re-submitting answers. As a final safeguard, our system automatically strips all answers and explanations from the data transmitted to the model, ensuring that only the question content is exposed and preventing answer leakage through data parsing.

2.3 Ranking System

To establish fair and stable rankings under dynamic evaluation conditions, we develop a calibrated ranking framework that combines LLM-as-a-Judge evaluation with relative scoring mechanisms. This approach ensures consistent and reliable model comparisons even when different question sets are used across evaluation sessions.

2.3.1 LLM-as-a-Judge Evaluation

To quantify answer quality, we establish a standardized scoring metric with an integer range of [0,3], ensuring consistent evaluation of response efficacy across diverse model architectures. For scoring implementation, we uniformly employ GPT-4o (OpenAI, 2023) as our judge, which has demonstrated high human-machine agreement through rigorous validation (detailed in Section 3.4). This choice of a single, validated scoring model eliminates potential biases introduced by varying evaluative criteria.

The scoring focuses on both core correctness and explanation quality, with core correctness serving as the primary indicator for score determination. The specific evaluation prompt and criteria are provided in Appendix B.

2.3.2 Evaluation Metrics

To mitigate systematic bias introduced by random sampling questions, LLMEval-Fair employs both relative score and absolute scores as evaluation metrics.

The absolute score S_{model} represents a model’s performance on $N = 1000$ questions, where each question receives a score s_i (with maximum score $s_{\text{max}} = 3$), mapped to the [0,100] interval:

$$S_{\text{model}} = \frac{\sum_{i=1}^N s_i}{N \times s_{\text{max}}} \times 100 \quad (1)$$

The relative score $R_{\text{SOTA}}^{\text{model}}$ is defined as the model’s absolute score relative to the current SOTA model’s absolute score on the same question set, mapped to the [0,100] interval:

$$R_{\text{SOTA}}^{\text{model}} = \frac{S_{\text{model}}}{S_{\text{SOTA}}} \times 100 \quad (2)$$

In our current evaluation, we use *Doubao-1.5-Thinking-Pro* as the reference SOTA model for relative score calculations.

3 Experiment and Analysis

Based on the three core challenges identified in our introduction, we design the LLMEval-Fair evaluation framework. In this section, we systematically investigate three critical research questions that arise from these challenges:

Research Question I: What authentic capability distributions and longitudinal trends do LLMs exhibit under LLMEval-Fair?

Research Question II: How does LLMEval-Fair dynamic evaluation compare with static benchmarks regarding ranking accuracy and contamination issues?

Research Question III: How stable and reliable is LLMEval-Fair’s relative ranking system under multi-round resampling and human-machine consistency validation?

Through comprehensive experiments designed to address these research questions, we aim to validate the effectiveness of our dynamic evaluation paradigm and provide empirical evidence for the superiority of contamination-resistant assessment frameworks.

3.1 Experimental Setup

3.1.1 Benchmarking LLMs on LLMEval-Fair

We tracked nearly 60 LLMs from June 2023 to December 2025, presenting full results in Appendix E

Model	$R_{\text{SOTA}}^{\text{model}}$	S_{model}	Eng.	Econ.	Edu.	Law	Lit.	Mgmt.	Sci.	Hist.	Med.	Mil.
<i>Open-source LLMs</i>												
DeepSeek-R1	97.40	91.23	9.47	9.43	9.27	9.37	8.83	9.37	9.03	9.53	8.50	8.43
DeepSeek-V3	96.47	90.36	9.30	9.57	8.93	9.23	8.60	9.13	8.97	9.47	8.83	8.33
Qwen-3-235B	96.42	90.32	9.23	9.43	9.03	9.50	8.23	9.43	8.97	9.17	8.73	8.60
Qwen-3-32B	92.22	86.38	8.43	9.10	8.57	9.10	7.77	9.47	8.67	9.30	7.70	8.27
<i>Closed-source LLMs</i>												
Doubao-1.5-Thinking-Pro	100.00	93.67	9.47	9.67	9.43	9.77	8.93	9.53	9.23	9.70	8.97	8.97
Gemini-2.5-Pro	97.22	91.07	9.20	9.47	9.20	9.30	8.43	9.63	9.07	9.40	8.50	8.87
Gemini-2.5-Pro-Thinking	97.15	91.00	9.13	9.50	9.37	9.47	8.40	9.63	9.20	9.27	8.30	8.73
Doubao-1.5-Pro	95.68	89.62	8.83	9.03	9.13	9.43	8.57	9.27	8.83	9.10	8.60	8.83
Kimi-K2	94.27	88.30	9.23	9.17	8.80	9.00	8.40	9.17	8.77	9.13	8.53	8.10
GPT-5	93.84	87.90	8.83	9.37	8.90	8.87	8.10	9.10	8.90	9.03	8.50	8.30
Claude-Sonnet-4.5-Thinking	93.48	87.57	8.90	9.17	8.80	8.97	8.00	9.23	8.90	9.00	8.27	8.33
o1	93.36	87.45	8.90	9.30	8.67	8.77	7.73	9.27	8.90	8.97	8.17	8.77
Claude-Sonnet-4-Thinking	91.03	85.27	8.57	9.00	8.63	8.73	7.57	9.10	8.93	8.70	7.97	8.07
Claude-Sonnet-4	91.00	85.24	8.57	8.80	8.50	8.70	7.80	9.03	8.80	8.80	8.17	8.07
GPT-4o-search	89.40	83.74	8.27	8.77	8.43	8.67	7.77	8.80	8.20	8.73	8.27	7.83
GPT-4o	88.09	82.51	7.90	8.67	8.30	8.33	7.17	8.97	8.57	8.67	7.63	8.30
o3-mini	84.13	78.80	7.97	8.60	8.30	8.20	6.73	8.57	8.53	7.17	7.03	7.70

Table 1: Overall and Subject-Level Scores. $R_{\text{SOTA}}^{\text{model}}$ represents the relative score (0-100 scale) as defined in Equation (2), with Doubao-1.5-Thinking-Pro as the reference SOTA model. S_{model} represents the absolute score (0-100 scale) as defined in Equation (1). Subject-level scores use a 10-point scale.

and focusing here on 17 representative models (proprietary and open-source). Appendix C provides a detailed overview of these representative models. Each model is evaluated across three prompting paradigms (Zero-Shot, Few-Shot, Chain-of-Thought) and 10 academic disciplines. In addition, we sample incorrect responses from all evaluated models and manually classify failures into five categories: disciplinary knowledge, misunderstanding, logical reasoning, factual inaccuracies, and format compliance.

3.1.2 Ablation Studies

We conduct comprehensive ablation studies across two key dimensions:

Benchmark Comparison: We measure Spearman correlation between LLMEval-Fair and static benchmarks (AGIEval (Zhong et al., 2024), C-Eval (Huang et al., 2023)) and perform fill-in-the-blank replay tests (1,000 questions, three attempts each) to assess contamination.

Ranking Validation: We conduct multi-round resampling (n=1000, 2000, 4000) to test stability. Human-machine agreement validation involves two independent rounds of human evaluation with Cohen’s κ coefficients computed against three LLM-as-Judge evaluators. Second, we run an ablation study comparing our relative ranking to the traditional Elo scoring system. The introduction to the Elo scoring system can be found in Appendix D.

3.2 Research Question I

Finding 1: All models converge to a performance ceiling of around 90% over a longitudinal period, with leading open-source LLMs rivaling proprietary SOTA. Figure 2 illustrates the performance growth trajectories of different model series over time. Table 1 presents comprehensive evaluation scores and subject-level breakdowns for each model under the LLMEval-Fair framework. Nearly all models approach a performance ceiling around 90 on academic knowledge tasks. This convergence indicates fundamental limits in current model architectures for knowledge-intensive evaluation.

Besides, the evaluation results show that top-performing models achieve remarkable scores: the proprietary Doubao-1.5-Thinking-Pro reaches 93.67 while the open-source DeepSeek-R1 achieves 91.23. These leading models substantially outperform established systems like GPT-4o at 82.51 and o3-mini. This demonstrates that with sufficient scale and training, both proprietary and open-source LLMs can achieve comparable performance.

Finding 2: Models demonstrate significant domain-specific performance variations, with specialized “thinking” abilities offering only marginal gains. Our analysis reveals a consistent pattern of domain-specific performance across all models. As shown in Table 1, models excel

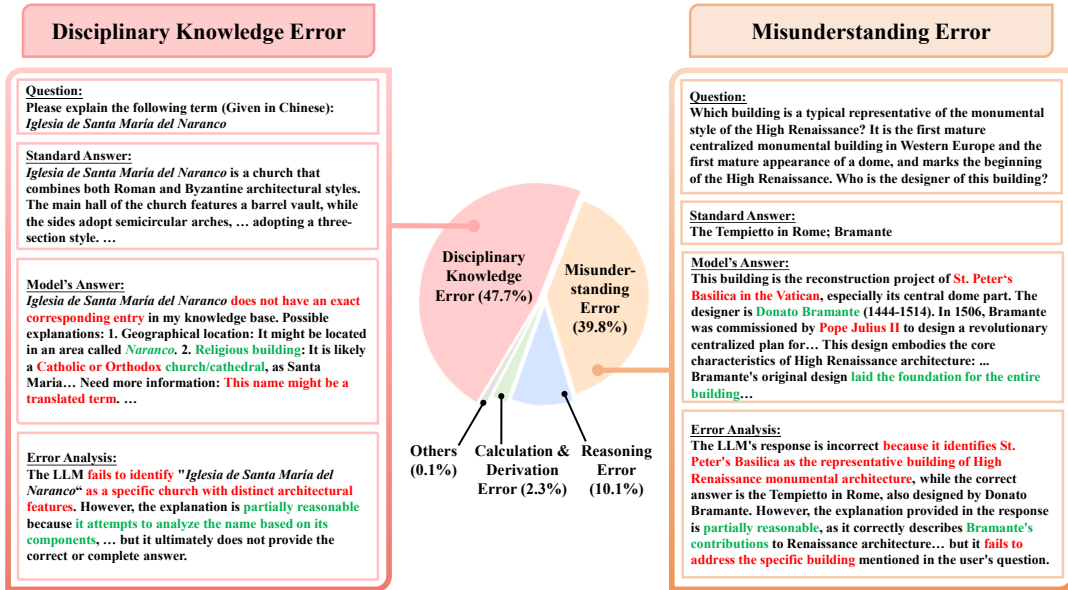


Figure 3: Distribution of model error causes and illustrative cases of the two most prevalent error types.

Severe Data Contamination. To assess the difference in leakage risk between publicly available static benchmarks and our private question bank, we conduct fill-in-the-blank replay tests on AGIEval (EN), AGIEval (ZH), C-Eval. For each dataset, we attempt 1000 questions with three completion attempts each, counting a question as successfully recalled if at least two attempts are correct. Table 4 reports the number of successful completions for each model.

We observe that on public static benchmarks almost all models achieve substantially higher completion counts, indicating extensive leakage of these questions into the training corpora. In contrast, on our private dataset, completion counts drop markedly, with most models near or below 100. Overall, these results suggest that public static benchmarks may suffer from significant leakage risk, whereas our private question bank, which is not included in large-scale pretraining corpora, supports a more contamination-resistant evaluation.

3.4 Research Question III

Finding 7: The relative ranking system demonstrates exceptional stability, with negligible variance across multi-round resampling and varying sample sizes. To validate the stability of the LLMEval-Fair ranking system, we conduct rigorous tests involving multi-round resampling and varying sample sizes. In our resampling test, we draw two independent question samples and find that the model ranking order remains identi-

Model	AGI (EN)	AGI (ZH)	C-Eval	Ours
DeepSeek-V3	97	153	136	80
ClaudeSonnet-4	179	248	224	179
Doubao-1.5	66	105	117	76
o3-mini	58	74	45	75
GPT-4o	54	86	62	48
Qwen3-32B	55	77	72	36

Table 4: Comparison of successful fill-in completions for different models on static benchmarks and LLMEval-Fair.

cal across both runs, with relative scoring exhibiting remarkably low variance (e.g., $\sigma^2 = 1.68$ for DeepSeek-V3 and just 0.25 for o3-mini). Furthermore, we test the system's sensitivity to question volume by running evaluations with sample sizes of $n=1000, 2000,$ and 4000 . Across all volumes, the ranking order holds firm. These results, detailed in Appendix F, provide strong empirical evidence that our relative ranking system is exceptionally robust and reliable.

Finding 8: The relative ranking system via LLM-as-Judge achieves high human-machine agreement and demonstrates superior robustness over alternative ranking systems. To validate our relative ranking system, we assess it on two critical dimensions: human alignment and methodological robustness. First, we measure its agreement with human experts by conducting two rounds of human evaluation and calculating Cohen's κ coefficient against our LLM-as-a-Judge. As

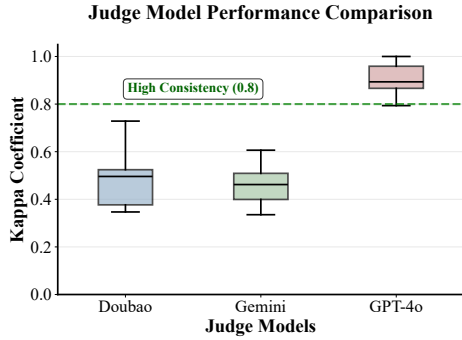


Figure 4: Cohen’s κ coefficients measuring agreement between human evaluators and three LLM judges across evaluations. GPT-4o achieves almost perfect agreement with human judgments ($\kappa = 0.907$).

shown in Figure 4, the system achieve near-perfect agreement with human judgments, confirming its reliability. Second, we conduct an ablation study comparing our relative ranking against an Elo-style rating baseline implemented under identical evaluation conditions. Specifically, both methods are evaluated on the same fixed pool of 1,000 queries, and the curves in Figure 5 are obtained by repeatedly sub-sampling from this pool to vary the effective sample size. The results show that our relative ranking consistently achieves higher Spearman’s ρ correlation and exhibits greater stability as the sample size changes. This dual validation provides strong evidence that our relative ranking system is both aligned with human evaluation and more robust than conventional rating-based alternatives.

4 Related Work

Static Knowledge-Intensive Benchmarks like MMLU (Hendrycks et al., 2021a) and C-Eval (Huang et al., 2023) measure factual knowledge and reasoning, yet their fixed question sets risk data leakage—inflating scores and misrepresenting capabilities (Banerjee et al., 2024; Xu et al., 2024a)—while promoting overfitting where models memorize test answers instead of generalizing (Deng et al., 2024b). This static paradigm, shared by LLMEval-1 and LLMEval-2 (Zhang et al., 2024), is now driving research toward more dynamic and robust evaluation methods.

Dynamic Human-Preference Evaluation (e.g., LMSYS (Zheng et al., 2024), Chatbot Arena (Chiang et al., 2024), AlpacaEval (Lab, 2023)) uses human preferences from pairwise or multi-turn interactions to assess conversational quality and real-time performance beyond static benchmarks. How-

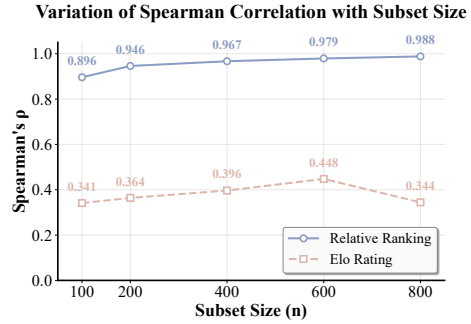


Figure 5: Relative ranking consistently outperforms Elo, reaching near-perfect correlation as the subset grows.

ever, these methods often lack depth in domain-specific reasoning, rely on non-expert annotators prone to bias in specialized assessments (Raju et al., 2024), and focus on dialogue at the expense of systematic academic knowledge coverage.

LLM-as-a-Judge employs language models as scalable evaluators of model outputs, as seen in G-Eval (Liu et al., 2023) and MT-Bench-101 (Bai et al., 2024). This paradigm offers high scalability, low cost, and strong human preference alignment. However, existing work focuses heavily on scoring mechanisms, leaving systematic judge calibration and bias mitigation underexplored (Zheng et al., 2023). Our work addresses this by rigorously analyzing and validating the paradigm through human-machine agreement studies.

5 Conclusion

We introduced LLMEval-Fair, a dynamic, contamination-resistant evaluation framework built on a private 220k-question bank, a two-layer anti-cheating architecture, and an LLM-as-Judge relative ranking pipeline. A 30-month longitudinal study of nearly 60 open-source and proprietary models revealed a consistent performance ceiling near 90%, systematic gaps in literature, medicine, and military knowledge, and widespread data leakage in static benchmarks. Our relative ranking method demonstrated negligible variance under multi-round resampling with varying sample sizes and achieved near-perfect agreement with human experts. We further confirmed that prompting format has minimal impact on performance in knowledge-intensive tasks, underscoring the superiority of dynamic, contamination-resistant evaluation over static benchmarks and the need for more trustworthy benchmarking practices.

566 Limitations

567 The main limitation is the sheer size of the ques-
568 tion bank, about 220,000 items, which makes a
569 truly comprehensive evaluation resource-intensive.
570 Running large-scale inference over the full set, val-
571 idating and aggregating results, and maintaining
572 the dataset through cleaning, deduplication, and
573 difficulty calibration require substantial compute,
574 time, and human effort. As a result, this evaluation
575 is more feasible for teams with stable infrastructure
576 and ongoing operational capacity, while smaller
577 groups may struggle to reproduce an equally com-
578 plete assessment with the same rigor and cadence.

579 References

580 Anthropic. 2025a. [Claude sonnet 4.5 system card](#). Sys-
581 tem card, Anthropic.

582 Anthropic. 2025b. [System card: Claude opus 4 &
583 claude sonnet 4](#). Technical Report / System Card.

584 Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jia-
585 heng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su,
586 Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024.
587 [Mt-bench-101: A fine-grained benchmark for evalu-
588 ating large language models in multi-turn dialogues](#).
589 In *Proceedings of the 62nd Annual Meeting of the
590 Association for Computational Linguistics (Volume
591 1: Long Papers)*, page 7421–7454. Association for
592 Computational Linguistics.

593 Sourav Banerjee, Ayushi Agarwal, and Eishkaran Singh.
594 2024. [The vulnerability of language model bench-
595 marks: Do they accurately reflect true LLM perfor-
596 mance?](#) *CoRR*, abs/2412.03597.

597 ByteDance. 2025. [Bytedance ai introduces doubao-
598 1.5-pro language model with a deep thinking mode](#).
599 MarkTechPost article.

600 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
601 Kaijie Zhu, Hao Chen, Linyi Yang, Xiaoyuan Yi,
602 Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang,
603 Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie.
604 2023. [A survey on evaluation of large language mod-
605 els](#). *CoRR*, abs/2307.03109.

606 Simin Chen, Pranav Pusarla, and Baishakhi Ray. 2025.
607 [Dynamic benchmarking of reasoning capabilities in
608 code large language models under data contamina-
609 tion](#). *CoRR*, abs/2503.04149.

610 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta-
611 sios Nikolas Angelopoulos, Tianle Li, Dacheng Li,
612 Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E.
613 Gonzalez, and Ion Stoica. 2024. [Chatbot arena: An
614 open platform for evaluating LLMs by human pref-
615 erence](#). In *Forty-first International Conference on
616 Machine Learning*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 617
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 618
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 619
Nakano, Christopher Hesse, and John Schulman. 620
2021. [Training verifiers to solve math word prob-
621 lems](#). *CoRR*, abs/2110.14168. 622

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, 623
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, 624
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, 625
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi- 626
hong Shao, Zhuoshu Li, Ziyi Gao, and 81 others. 627
2025. [Deepseek-r1: Incentivizing reasoning capa-
628 bility in llms via reinforcement learning](#). *CoRR*,
629 abs/2501.12948. 630

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx- 631
uan Wang, Bochao Wu, Chengda Lu, Chenggang 632
Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, 633
Damai Dai, Daya Guo, Dejian Yang, Deli Chen, 634
Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, 635
and 81 others. 2024. [Deepseek-v3 technical report](#).
636 *CoRR*, abs/2412.19437. 637

Jasper Dekoninck, Mark Niklas Müller, Maximilian 638
Baader, Marc Fischer, and Martin T. Vechev. 2024. 639
[Evading data contamination detection for language
640 models is \(too\) easy](#). *CoRR*, abs/2402.02823. 641

Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Ger- 642
stein, and Arman Cohan. 2024a. [Investigating data
643 contamination in modern benchmarks for large lan-
644 guage models](#). In *Proceedings of the 2024 Con-
645 ference of the North American Chapter of the As-
646 sociation for Computational Linguistics: Human
647 Language Technologies (Volume 1: Long Papers)*,
648 *NAACL 2024, Mexico City, Mexico, June 16-21, 2024*,
649 pages 8706–8719. Association for Computational
650 Linguistics. 651

Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Ger- 652
stein, and Arman Cohan. 2024b. [Investigating data
653 contamination in modern benchmarks for large lan-
654 guage models](#). In *Proceedings of the 2024 Con-
655 ference of the North American Chapter of the Associ-
656 ation for Computational Linguistics: Human Lan-
657 guage Technologies (Volume 1: Long Papers)*, pages
658 8706–8719, Mexico City, Mexico. Association for
659 Computational Linguistics. 660

Gemini Team, Google DeepMind. 2025. [Gemini 2.5:
661 Pushing the frontier with advanced reasoning, multi-
662 modality, long context, and next generation agentic
663 capabilities](#). Technical report, Google DeepMind. 664

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, 665
Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 666
2021a. [Measuring massive multitask language under-
667 standing](#). In *International Conference on Learning
668 Representations*. 669

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul 670
Arora, Steven Basart, Eric Tang, Dawn Song, and
671 Jacob Steinhardt. 2021b. [Measuring mathematical
672](#)

673	problem solving with the MATH dataset. In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	728
674		729
675		730
676		731
677		
678	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, and 1 others. 2023. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. <i>Advances in Neural Information Processing Systems</i> , 36:62991–63010.	732
679		733
680		734
681		
682		
683		
684		
685	Michael B. Jones, John Bradley, and Nat Sakimura. 2015. <i>JSON Web Token (JWT)</i> . RFC 7519.	
686		
687	Tatsu Lab. 2023. AlpacaEval: An automatic evaluator for instruction-following language models. https://github.com/tatsu-lab/alpaca-eval . Accessed: 2025-07-31.	
688		
689		
690		
691	Md. Tahmid Rahman Laskar, Sawsan Alqahtani, M. Saiful Bari, Mizanur Rahman, Mohammad Abdullah Matin Khan, Haidar Khan, Israt Jahan, Amran Bhuiyan, Chee-Wei Tan, Md. Rizwan Parvez, Enamul Hoque, Shafiq Joty, and Jimmy Huang. 2024. A systematic survey and critical review on evaluating large language models: Challenges, limitations, and recommendations. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 13785–13816. Association for Computational Linguistics.	
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
703	Songyang Liu, Chaozhuo Li, Jiameng Qiu, Xi Zhang, Feiran Huang, Litian Zhang, Yiming Hei, and Philip S. Yu. 2025. The scales of justitia: A comprehensive survey on safety evaluation of llms. <i>CoRR</i> , abs/2506.11094.	
704		
705		
706		
707		
708	Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	
709		
710		
711		
712		
713		
714		
715	Moonshot AI. 2025. <i>Kimi k2 technical report</i> . Technical report, Moonshot AI.	
716		
717	OpenAI. 2023. <i>GPT4 technical report</i> .	
718	OpenAI. 2024. <i>o1 system card</i> . Technical report, OpenAI.	
719		
720	OpenAI. 2025a. <i>Gpt-5 system card</i> . System card, OpenAI.	
721		
722	OpenAI. 2025b. <i>o3-mini system card</i> . Technical report, OpenAI.	
723		
724	Ravi Raju, Swayambhoo Jain, Bo Li, Jonathan Li, and Urmish Thakker. 2024. <i>Constructing domain-specific evaluation sets for llm-as-a-judge</i> . <i>Preprint</i> , arXiv:2408.08808.	
725		
726		
727		
	Cheng Xu, Shuhao Guan, Derek Greene, and M-Tahar Kechadi. 2024a. <i>Benchmark data contamination of large language models: A survey</i> . <i>Preprint</i> , arXiv:2406.04244.	728
		729
		730
		731
	Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. 2024b. <i>Benchmarking benchmark leakage in large language models</i> . <i>Preprint</i> , arXiv:2404.18824.	732
		733
		734
	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, and 41 others. 2025. <i>Qwen3 technical report</i> . <i>CoRR</i> , abs/2505.09388.	735
		736
		737
		738
		739
		740
		741
	Yue Zhang, Ming Zhang, Haipeng Yuan, Shichun Liu, Yongyao Shi, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. LlmEval: A preliminary study on how to evaluate large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19615–19622.	742
		743
		744
		745
		746
		747
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2024. <i>LMSYS-chat-1m: A large-scale real-world LLM conversation dataset</i> . In <i>The Twelfth International Conference on Learning Representations</i> .	748
		749
		750
		751
		752
		753
		754
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. <i>Judging LLM-as-a-judge with MT-bench and chatbot arena</i> . In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	755
		756
		757
		758
		759
		760
		761
	Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2024. <i>Agieval: A human-centric benchmark for evaluating foundation models</i> . In <i>Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024</i> , pages 2299–2314. Association for Computational Linguistics.	762
		763
		764
		765
		766
		767
		768
		769
	A Dataset	770
	This section provides supplementary information on our LLMEval-Fair dataset employed in the study, to clarify the academic disciplines and question types covered, and to elaborate methodologies for scaling the questions.	771
		772
		773
		774
		775
	A.1 Categories of Academic Disciplines	776
	A.1.1 Categories of Academic Disciplines	777
	As illustrated in Figure 6, we collected graduate-level examination questions spanning 13 primary (Philosophical Sciences, Economic Sciences, Law,	778
		779
		780

Education, Literature, History, Engineering, Agronomy, Medicine, Military Science, Management Sciences, Arts, and Sciences) and more than 50 secondary academic disciplines recognized by China’s Ministry of Education. Two-thirds of the questions are derived from Chinese universities’ Postgraduate Entrance Exams, and one-third are from Undergraduate Final Exams of comparable difficulty. Detailed distribution of question sources is listed in Table 5.

A.1.2 Categories of Question Types

The raw dataset encompasses a diverse range of original question formats, including multiple-choices, fill-in-the-blank, true-or-false, short-answer, term explanation, and material analysis questions. Following question expansion and formatting, we have unified all questions and their answers into a fill-in-the-blank-like question-answer format, abandoning the complex answer structures of various question types, such as options A to E in multiple-choice questions, and “true” or “false” in true-or-false questions. This natural question-answer format enables the dataset to more thoroughly showcase the model’s capabilities.

A.2 Details of Expanding the Dataset

As shown in Table 6, we have amassed a substantially large quantity of original questions, with a marked surge in numbers following the expansion of our dataset.

A.2.1 Original Data Construction Pipeline

Original questions, structured simply and organized by subject for initial collation, follow this construction pipeline: first, converting Excel, Word, and PDF test papers to TXT; then batch splitting into JSON-formatted questions via scripts; and finally conducting data screening.

The latter involves three steps: (1) Expert review removes factually erroneous or irrelevant questions. (2) Batch splitting classifies and isolates questions, addressing errors such as content overlap, missing questions, or misclassifications. (3) Format cleaning resolves encoding conflicts, special characters, symbol consistency, redundancy, and typos.

A.2.2 Data Expanding Pipeline

To accommodate diverse application scenarios, this study proposes an augmented data format that complements the original question structure. The augmented dataset incorporates comprehensive metadata, including primary disciplinary cate-

gories, secondary disciplinary categories, question descriptions, answer content, and unique identifiers (UUIDs), which facilitates categorized data management. An example of expanding the Multiple-Choices question is shown in Figure 7.

Two additional key verification procedures are implemented upon the core augmentation strategy: (1) Format validation, which entails checking the consistency of option counts for multiple-choice questions and the alignment of answer spaces for Fill-in-the-Blank questions; and (2) Redundancy checks, which involve detecting duplicates among split questions and ensuring the uniqueness of question UUIDs. The data formats before and after expanding is illustrated in Figure 8.

B Prompts

This section presents the complete set of prompts used in LLMEval-Fair for different evaluation paradigms. We provide the specific prompt templates for few-shot learning, chain-of-thought reasoning, and LLM-based automated evaluation to ensure reproducibility of our experimental results.

The prompts used for testing few-shot and chain-of-thought methods are shown in Figures 9 and 10, respectively. The prompt used for LLM-based evaluation is shown in Figure 11.

C Evaluation Model Selection

To ensure a comprehensive and rigorous assessment, we conducted a preliminary evaluation on a broad set of 59 large language models. From this extensive pool, we selected a representative subset of 17 models for the detailed analysis presented in the main text. The detailed release dates for these selected models are provided in Table 7.

The proprietary frontier models include the **OpenAI** lineup, featuring the widely deployed **GPT-4o** and **GPT-4o-search** (OpenAI, 2023), alongside the next-generation flagship **GPT-5** (OpenAI, 2025a). Furthermore, we assess the **o1** series (OpenAI, 2024) and the efficient **o3-mini** (OpenAI, 2025b), which represent specialized reasoning models trained with large-scale reinforcement learning to solve complex scientific and mathematical problems.

We include models capable of extended reasoning from other major providers. **Anthropic** is represented by the standard **Claude-Sonnet-4** (Anthropic, 2025b), as well as the **Claude-Sonnet-4-Thinking** and the advanced **4.5-Thinking** variants

Type	Number of Topics	Proportion (%)
Undergraduate Final Exams	26633	34.1
Postgraduate Entrance Exams	51376	65.9
Total	78009	100.0

Undergraduate Final Exams	71038	31.1
Postgraduate Entrance Exams	157566	68.9
Total	228604	100.0

Table 5: Distribution of question number and proportions for Undergraduate Final Exams and Postgraduate Entrance Exams.

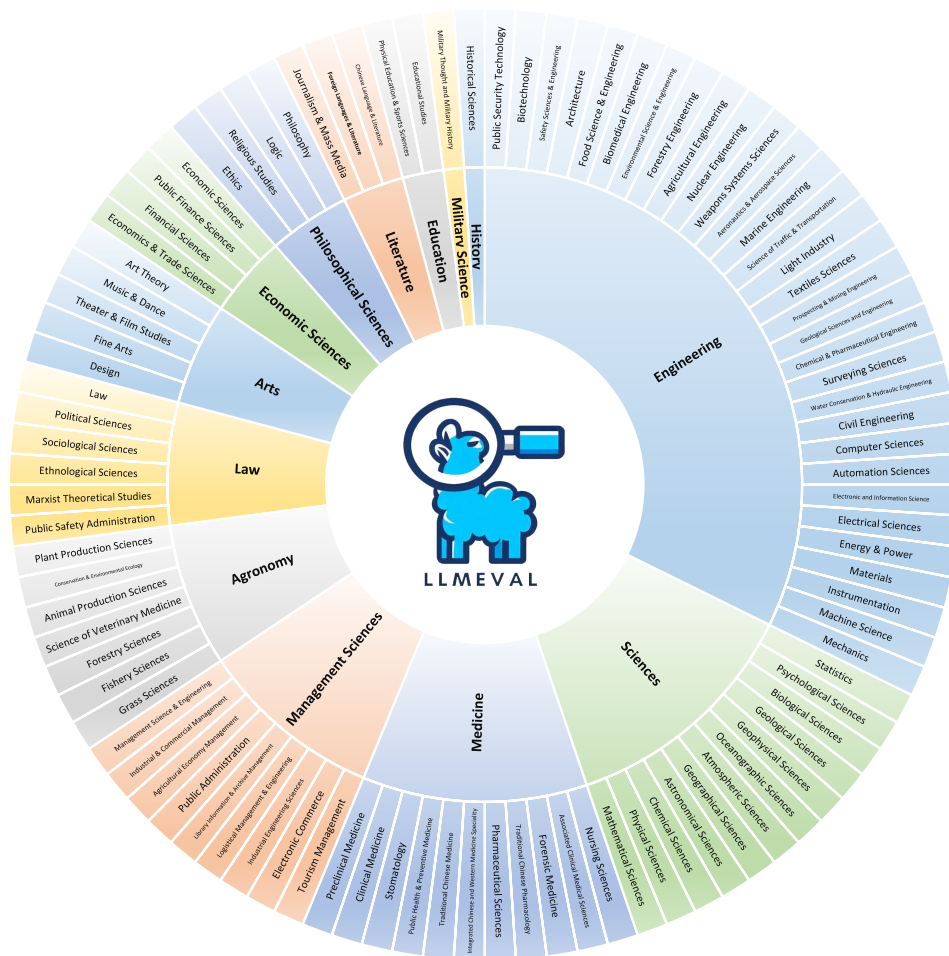


Figure 6: Categories of Primary and Secondary Academic Disciplines.

Subject	Original Count	Rewritten Count	Increase Count	Increase Percentage (%)
Philosophical Sciences	2194	10969	8775	399.95
Medicine	30772	109974	79202	257.38
Law	9262	30116	20854	225.16
Management Sciences	2448	7945	5497	224.55
Engineering	4926	13263	8337	169.24
Sciences	6182	15669	9487	153.46
Economic Sciences	9245	18124	8879	96.04
Military Science	611	1187	576	94.27
Education	2781	5094	2313	83.17
History	1749	3178	1429	81.70
Literature	7839	13085	5246	66.92
Total	78009	228604	150595	193.05

Table 6: Distribution of Original and Rewritten Counts Across Disciplines.

(Anthropic, 2025a), which operate in a specialized mode to perform self-reflection before generating responses. Similarly, Google’s contribution consists of Gemini-2.5-Pro (Gemini Team, Google DeepMind, 2025) and its reasoning-enhanced counterpart, Gemini-2.5-Pro-Thinking (Gemini Team, Google DeepMind, 2025), which incorporates internal chain-of-thought processes.

Regarding open-weights architectures and diverse scaling strategies, the study includes the DeepSeek family: DeepSeek-V3 (DeepSeek-AI et al., 2024), a strong Mixture-of-Experts (MoE) model, and DeepSeek-R1 (DeepSeek-AI et al., 2025), a model specifically optimized for reasoning tasks through post-training. The Qwen-3 series is also evaluated to represent distinct points on the parameter spectrum, featuring both the massive 235B model and the compact 32B variant (Yang et al., 2025).

Finally, we evaluate other high-performing systems to ensure a broad representation of the current landscape. This includes Moonshot’s Kimi-K2 (Moonshot AI, 2025) and the Doubao-1.5 series. Notably, Doubao-1.5-Thinking-Pro (ByteDance, 2025) serves as a key reference in our study, having demonstrated state-of-the-art capabilities in preliminary screenings.

D The Elo Rating System

The Elo rating system is a widely recognized method for calculating relative skill levels in zero-sum games, originally developed for chess and recently popularized for evaluating Large Language Models (e.g., Chatbot Arena). This system derives ratings from the outcomes of pairwise com-

Table 7: The representative list of models selected from a total of 59 evaluated LLMs. Dates for unreleased models are estimated based on technical previews.

Model	Released Date
Doubao-1.5-Thinking-Pro (2025)	April 15, 2025
Doubao-1.5-Pro (2025)	January 15, 2025
Gemini-2.5-Pro-Thinking (2025)	June 5, 2025
Gemini-2.5-Pro (2025)	June 5, 2025
DeepSeek-R1 (2025)	May 28, 2025
DeepSeek-V3 (2024)	March 24, 2024
Qwen-3-235B (2025)	April 29, 2025
Qwen-3-32B (2025)	April 29, 2025
GPT-5 (2025a)	June 7, 2025
Claude-Sonnet-4.5-Thinking (2025a)	September 29, 2025
Claude-Sonnet-4-Thinking (2025b)	May 14, 2025
Claude-Sonnet-4 (2025b)	May 14, 2025
o1 (2024)	December 17, 2024
o3-mini (2025b)	January 29, 2025
GPT-4o-search (2023)	November 20, 2024
GPT-4o (2023)	November 20, 2024
Kimi-K2 (2025)	September 5, 2025

parisons, providing a probabilistic framework to predict the likelihood of one entity outperforming another. Given two entities A and B with current ratings R_A and R_B , the expected score E_A (representing the probability of A winning) is calculated using a logistic curve:

$$E_A = \frac{1}{1 + 10^{(R_B - R_A)/400}} \quad (3)$$

Following a match, the ratings are updated based on the discrepancy between the actual outcome S_A (where 1 represents a win, 0 a loss, and 0.5 a tie) and the expected probability. The update rule is given by $R'_A = R_A + K(S_A - E_A)$, where K is

Original Multiple-Choices questions:

Title: It is known that the sixth level of a complete binary tree (let the root be the first level) has eight leaves, then the number of nodes of the complete binary tree is at most ().

- A. 39
- B. 52
- C. 111
- D. 119

Answer: C

Expanded questions:

Question 1: It is known that the sixth level of a complete binary tree (let the root be the first level) has eight leaves, then the number of nodes of the complete binary tree is at most: ()

Is it correct to place the answer "39" in the provided space?

Answer: False

Question 2: It is known that the sixth level of a complete binary tree (let the root be the first level) has eight leaves, then the number of nodes of the complete binary tree is at most: ()

Is it correct to place the answer "52" in the provided space?

Answer: False

Question 3: It is known that the sixth level of a complete binary tree (let the root be the first level) has eight leaves, then the number of nodes of the complete binary tree is at most: ()

Is it correct to place the answer "111" in the provided space?

Answer: True

Question 4: It is known that the sixth level of a complete binary tree (let the root be the first level) has eight leaves, then the number of nodes of the complete binary tree is at most: ()

Is it correct to place the answer "119" in the provided space?

Answer: False

Figure 7: Example of how to expand a Multiple Choice question.

925 a constant factor that determines the sensitivity of
926 the rating adjustment.

927 **E LLMEval-Fair Leaderboard**

928 This section presents comprehensive evaluation results
929 from our longitudinal study tracking over 50
930 LLMs from late 2023 to mid-2025. We provide
931 complete performance rankings and analyze the
932 consistency of model capabilities across different
933 prompting paradigms.

934 We tracked over 50 LLMs from late 2023 to mid-
935 2025. Here, we present the complete evaluation
936 results of our model assessments. The comprehen-
937 sive results, including scores for all models across

10 academic disciplines, are presented in Table 8.

The models we selected in main paper was evalu-
ated across three prompting paradigms: Zero-Shot
(ZS), Few-Shot (FS), and Chain-of-Thought (CoT).
As shown in Table 2, the performance variance
across these paradigms remains below 1.6 points
for all evaluated models, indicating that core ca-
pabilities are not significantly influenced by the
prompting format.

F Details in Experiment

This section provides additional experimental de-
tails demonstrating the robustness and stability of
our ranking system. We present validation experi-

938
939
940
941
942
943
944
945
946
947
948
949
950

```

category: Engineering,
sub_category: Computer Science,
question_uuid: a4b92619-f0ac-xxxx-xxxx-9f11ce3b4704,
dialog:
role: user,
content: Judge the correctness of the following statements, answer true/false, and give your
reasons. The continue statement in a loop breaks out of the current loop.
hint: false,

gpt4res: Answer: false. The reason is that the continue statement does not break out of the entire
loop.It skips the rest of the current loop and immediately moves on to the next iteration of the loop
instead of ending the current loop completely.
gpt4judge:
judge: Overall Rating: 3 The reason why you gave this rating: The LLM provided the correct
answer (false) and gave a reasonable explanation. It accurately explained that the ‘continue‘
statement does not terminate the entire loop but instead skips the remaining part of the current
iteration and proceeds to the next iteration of the loop. This explanation aligns with the correct
behavior of the ‘continue‘ statement in programming.,
stars: 3,
reason: The LLM provided the correct answer (false) and gave a reasonable explanation.
It accurately explained that the ‘continue‘ statement does not terminate the entire loop but
instead skips the remaining part of the current iteration and proceeds to the next iteration of the
loop. This explanation aligns with the correct behavior of the ‘continue‘ statement in programming.

```

Figure 8: Detailed entries of a single question after expanding.

ments across different sample sizes and show the consistency of our relative scoring methodology.

F.1 Sampling Validation

To verify the stability of our ranking system, we conducted evaluations across multiple sample sizes of $n=1000$ (in three rounds), 2000, and 4000 questions. The results in Table 9 demonstrate remarkable consistency in ranking order across all sample sizes.

Our relative scoring methodology produces smaller variance compared with absolute scoring approaches. The variance analysis reveals that top tier models show exceptional stability. DeepSeek-V3 exhibits a variance of 0.51 and O3-mini exhibits a variance of 0.95 while GPT-4o exhibits the highest variance of 1.63. Even the maximum variance represents less than two percent fluctuation indicating robust measurement precision.

The three independent one-thousand-sample runs demonstrate high reproducibility with models maintaining consistent relative positions across all test conditions. These findings validate that our

ranking methodology captures stable model capabilities rather than random fluctuations.

F.2 LLM-as-Judge Validation

We calculated Cohen’s κ coefficients between human evaluations and three LLM judges Doubao Gemini and GPT-4o across two evaluation rounds for thirteen models. As shown in Table 10 GPT-4o demonstrates superior performance with κ values consistently above 0.90 with an average of 0.901 in Round 1 and 0.892 in Round 2 indicating almost perfect agreement.

The data reveal notable stability differences among judges. GPT-4o maintains consistently high agreement across both rounds with minimal variation while Doubao and Gemini exhibit more fluctuation between rounds. Specifically Doubao’s performance ranges from 0.232 to 0.745 across different models and Gemini exhibits even greater instability with some models showing dramatic drops between rounds—for example DeepSeek-V3 declines from 0.627 to 0.147.

In contrast Doubao and Gemini show lower over-

995 all agreement with average κ values of 0.493 and
996 0.446 for Round 1 and Round 2 for Doubao and
997 0.494 and 0.400 for Gemini. Based on GPT-4o’s
998 consistently high correlation with human evalua-
999 tions and superior stability we selected it as our
1000 primary judge for reliable assessment.

1001 **G Implementation Details**

1002 This section provides detailed information about
1003 the annotation processes and evaluation procedures
1004 underlying our LLMEval-Fair platform. We de-
1005 scribe the expert involvement in data curation, vali-
1006 dation processes, and the associated costs to ensure
1007 transparency and reproducibility.

1008 **G.1 Data Annotation Process**

1009 A total of 38 experts were engaged in data anno-
1010 tation and cleaning processes, with an average of
1011 more than 3 relevant specialists assigned to each
1012 discipline. For the annotation of original data, to
1013 mitigate fatigue-induced errors, annotation tasks
1014 for each expert were distributed across a 30–60 day
1015 period.

1016 The cumulative remuneration disbursed to ex-
1017 perts involved in data annotation and cleaning
1018 amounted to \$48,700. Ongoing investments are
1019 being allocated to further hire experts to expand
1020 the dataset.

1021 **G.2 Manual Evaluation Process**

1022 To validate our LLM-as-Judge approach, we con-
1023 ducted comprehensive human evaluation studies.
1024 A total of 18 experts participated in manual eval-
1025 uation processes, with an average of about 2 rele-
1026 vant specialists assigned to each discipline. This
1027 expert-based validation ensures that our automated
1028 evaluation system maintains high agreement with
1029 human judgment standards.

1030 The evaluation process involved multiple rounds
1031 of assessment across 13 representative models,
1032 with experts providing independent judgments that
1033 were subsequently compared against our LLM-
1034 based evaluation system using Cohen’s κ coeffi-
1035 cient.

1036 **G.3 Cost Analysis**

1037 The development and validation of LLMEval-Fair
1038 required substantial investment in both human ex-
1039 pertise and computational resources. We spent
1040 more than \$5,000 on using latest APIs of LLMs and
1041 deploying models for evaluation purposes. Addi-
1042 tionally, \$10,000 was allocated for hiring qualified

volunteers to conduct manual evaluations, ensuring
rigorous validation of our automated assessment
framework.

1043 **H JWT Authentication Process**

1044 This section describes the detailed implementation
1045 of our JSON Web Token (JWT) authentication sys-
1046 tem, which forms the outer layer of our two-tier
1047 anti-cheating architecture. The JWT process en-
1048 sures secure and authenticated access to our evalua-
1049 tion platform while preventing unauthorized access
1050 and session manipulation.
1051
1052
1053

1054 Our JWT implementation follows a standard
1055 three-phase protocol: token generation, transmis-
1056 sion, and verification. Algorithm 1 outlines the
1057 complete JWT authentication workflow used in
1058 LLMEval-Fair.

Algorithm 1: JWT Authentication Process in LLMEval-Fair

- 1: **Server-side Token Generation:**
- 2: Generate a unique user identity (`user_id`)
- 3: Generate current timestamp and expiration time (`exp`)
- 4: Construct payload \leftarrow {`user_id`, timestamp, `exp`, `session_id`, `permissions`}
- 5: Sign payload with server Secret using HMAC-SHA256 to generate JWT
- 6: **return** JWT to the authenticated user
- 7:
- 8: **Client-side Request:**
- 9: Include JWT in Authorization header: Bearer <token>
- 10: Send request to evaluation endpoint
- 11:
- 12: **Server-side Verification:**
- 13: Extract JWT from Authorization header
- 14: Verify JWT signature using server Secret
- 15: Parse payload and extract claims
- 16: **if** JWT signature is invalid **then**
- 17: **return** HTTP 401 Unauthorized
- 18: **end if**
- 19: **if** `current_time` > `exp` **then**
- 20: **return** HTTP 401 Token Expired
- 21: **end if**
- 22: **if** `user_id` not found **or** `permissions` insufficient **then**
- 23: **return** HTTP 403 Forbidden
- 24: **end if**
- 25: **if** session validation fails (e.g., concurrent sessions detected) **then**
- 26: **return** HTTP 403 Session Invalid
- 27: **end if**
- 28: Allow evaluation operation to proceed
- 29: Log access attempt with `user_id`, timestamp, and `session_id`

Input:

Here are several examples:

Question:

Please determine the correctness of the following statement. Answer with true/false and provide a reason. If all elements below the diagonal in the adjacency matrix of a directed graph are zero, then the graph must have a topological ordering.

Answer:

true. Because if all elements below the diagonal in the adjacency matrix of a directed graph are zero, the graph is a Directed Acyclic Graph (DAG), so it must have a topological ordering.

Question:

Please explain the following term: "Double Hundred Policy".

Answer:

It refers to "let a hundred flowers bloom, let a hundred schools of thought contend." This was a policy on science and culture officially proposed by Mao Zedong in 1956 and confirmed by the Central Committee of the Communist Party of China. The policy was severely undermined after 1957, but was reestablished and implemented following the Third Plenary Session of the 11th Central Committee.

Question:

Please determine the correctness of the following statement. Answer with true/false and provide a reason. According to the convertibility theory, the scope of commercial banks' assets expanded from short-term turnover loans to consumer loans.

Answer:

false. Convertibility Theory: also known as the asset conversion theory. This theory suggests that to maintain liquidity for withdrawals, commercial banks can invest part of their funds in transferable securities. Since these profitable assets can be sold at any time and converted into cash, loans are not necessarily limited to short-term and self-liquidating types. Clearly, this theory emerged in the context of developing financial instruments and markets. Significance: it expanded the scope of bank asset operations. Drawbacks: it does not guarantee that assets can be liquidated without capital loss (which requires high asset quality and stable market conditions); it is also constrained by central bank monetary policy (e.g., the risk of a rise in discount rates). This theory provides a theoretical basis for Chinese commercial banks to engage in securities business (investment operations). However, in China, commercial banks' investment activities are restricted due to: 1. limited investment instruments and underdeveloped credit mechanisms; 2. management systems narrowing investment scopes, and separation of operations preventing commercial banks from investing in stocks; 3. the nature of state-owned commercial banks limits their willingness for autonomous investment.

Question:

Why can the results of animal experiments not be fully applied to clinical practice?

Answer:

Because there are differences between humans and animals not only in cellular morphology and metabolism, but also due to the highly developed human nervous system, which is associated with language and thought (the second signaling system). Although there are similarities, the essential differences mean that human diseases cannot all be replicated in animals. Even if they can be replicated, animal responses are simpler than human responses. Therefore, results from animal experiments cannot be mechanically and fully applied to clinical practice without analysis. Only by comparing, analyzing, and synthesizing animal experiment results with clinical data can they be used as references in clinical medicine and provide a basis for studying the causes, mechanisms, prevention, and treatment of clinical diseases.

The following is the question to be answered:

Question:{question}

Answer:

Figure 9: The Few-shot Prompt Template.

Input:

{question}

Please think step by step and provide the final answer.

Figure 10: The Chain-of-Thought Prompt Template.

Input:

Please evaluate the following response from the LLM regarding a discipline-specific question based on the following criteria. You must score it on a scale of 0, 1, 2 or 3 stars:

Overall Rating:

0 star indicates wrong answer with a wrong explanation

1 stars indicate wrong answer but a partially reasonable explanation

2 stars indicate a correct answer with a partially reasonable explanation

3 stars indicate an correct answer with a reasonable explanation

User: {question}

LLM: {LLM response}

The correct answer to user's question is: {correct answer}

You must provide your feedback in the following format:

"Overall Rating":numbers of its stars(int)

The reason why you gave this rating: <Your Reason>(str)

Figure 11: The Prompt Template for LLM Judgement.

Model	R_{SOTA}^{model}	S_{model}	Eng.	Econ.	Edu.	Law	Lit.	Mgmt.	Sci.	Hist.	Med.	Mil.
<i>Open-source LLMs</i>												
DeepSeek-R1	97.40	91.23	9.47	9.43	9.27	9.37	8.83	9.37	9.03	9.53	8.50	8.43
DeepSeek-V3	96.47	90.36	9.30	9.57	8.93	9.23	8.60	9.13	8.97	9.47	8.83	8.33
Qwen3-235B	96.42	90.32	9.23	9.43	9.03	9.50	8.23	9.43	8.97	9.17	8.73	8.60
QwQ-32B	94.51	88.53	8.30	9.46	9.23	9.33	7.83	9.46	8.65	9.27	8.57	8.43
DeepSeek-V3.2	92.27	86.43	8.73	9.13	8.53	8.70	7.40	9.33	8.87	9.37	8.53	7.83
Qwen3-32B	92.22	86.38	8.43	9.10	8.57	9.10	7.77	9.47	8.67	9.30	7.70	8.27
GLM-4-32B	88.43	82.83	7.77	8.97	8.33	8.33	7.03	9.13	8.27	8.77	8.23	8.00
Qwen2.5-32B-Instruct	85.06	79.68	7.70	8.57	8.33	8.33	6.70	8.50	8.17	7.70	7.60	8.08
Qwen-Turbo-1101	83.72	78.42	7.97	8.37	8.03	8.23	6.40	8.50	8.10	7.50	7.27	8.05
Yi-34B-Chat	70.15	65.71	5.77	6.63	7.37	7.53	5.47	5.77	5.47	7.47	6.30	7.93
Megrez-3B-Instruct	67.01	62.77	5.80	6.77	6.80	7.13	5.40	6.87	5.70	6.53	5.70	6.07
Qwen2-7B-Instruct	65.15	61.03	5.47	6.73	6.33	7.60	5.13	6.17	6.17	5.73	5.33	6.37
Nanbeige-Plus	65.10	60.98	5.78	5.57	6.77	7.37	5.37	5.93	5.45	6.30	5.67	6.77
Phi-4-Final	63.98	59.93	5.80	6.47	6.23	6.53	5.53	6.30	6.27	5.50	5.43	5.87
Llama-3.2-90B-Vision	61.74	57.83	5.63	6.33	6.20	5.80	4.73	6.10	6.57	5.03	5.27	6.17
Llama-3.3-70B	60.85	57.00	5.80	6.90	5.63	5.70	5.47	5.70	6.30	4.70	4.87	5.93
Baichuan2-13B-Chat	58.28	54.59	4.47	5.53	7.40	6.90	4.63	4.80	4.33	6.23	4.60	5.70
Qwen-plus	56.58	53.00	4.40	5.10	6.53	6.53	5.00	4.77	4.87	5.17	5.13	5.50
Qwen-turbo	55.76	52.23	4.10	6.07	6.63	6.43	4.43	4.53	4.97	5.27	4.37	5.43
Nanbeige-16B	55.45	51.94	4.37	5.30	6.50	6.30	3.97	4.70	4.07	5.90	4.73	6.10
Mixtral-8x7B-Instruct	51.69	48.42	4.27	5.47	6.47	6.40	3.13	4.50	5.07	3.57	4.37	5.17
ChatGLM-2-6B	42.31	39.63	2.33	3.77	5.97	6.13	2.83	3.83	2.60	3.80	4.00	4.37
Llama-3.1-8B	41.25	38.64	3.87	4.20	4.27	4.17	3.50	3.83	4.30	3.17	3.20	4.13
Ziya-13B-v1.1	40.18	37.64	2.77	3.97	5.17	5.33	2.80	3.77	2.53	3.70	3.03	4.57
InternLM-7B-Chat	38.71	36.26	2.63	3.67	4.87	5.57	3.17	3.33	2.33	4.03	3.13	3.53
Linly-LLaMA2-13B	37.03	34.69	2.20	3.77	4.50	5.00	2.43	3.33	2.53	3.90	2.50	4.53
Phi-3-Medium-128K	36.95	34.61	2.27	4.17	3.70	4.23	2.87	4.50	3.57	3.20	2.27	3.83
BELLE-Llama2-13B-Chat	36.25	33.96	2.57	3.07	4.93	4.73	2.83	3.80	2.43	3.33	2.40	3.87
Llama-2-7B-Chat	25.22	23.62	1.53	3.43	3.00	3.73	1.73	2.43	1.97	2.17	0.80	2.83
<i>Closed-source LLMs</i>												
Doubao-1.5-Thinking-Pro	100.00	93.67	9.47	9.67	9.43	9.77	8.93	9.53	9.23	9.70	8.97	8.97
Gemini-2.5-Pro	97.22	91.07	9.20	9.47	9.20	9.30	8.43	9.63	9.07	9.40	8.50	8.87
Gemini-2.5-Pro-Thinking	97.15	91.00	9.13	9.50	9.37	9.47	8.40	9.63	9.20	9.27	8.30	8.73
Doubao-1.5-Pro	95.68	89.62	8.83	9.03	9.13	9.43	8.57	9.27	8.83	9.10	8.60	8.83
GLM-4.6	95.26	89.23	8.80	9.27	8.70	9.23	8.40	9.63	8.90	9.30	8.43	8.57
Kimi-K2	94.27	88.30	9.23	9.17	8.80	9.00	8.40	9.17	8.77	9.13	8.53	8.10
GPT-5	93.84	87.90	8.83	9.37	8.90	8.87	8.10	9.10	8.90	9.03	8.50	8.30
Claude-Sonnet-4.5-Thinking	93.48	87.57	8.90	9.17	8.80	8.97	8.00	9.23	8.90	9.00	8.27	8.33
o1	93.36	87.45	8.90	9.30	8.67	8.77	7.73	9.27	8.90	8.97	8.17	8.77
Claude-Sonnet-4.5	93.31	87.40	8.80	8.97	8.93	8.73	8.37	9.10	8.97	8.93	8.13	8.47
Gemini-2.5-Flash-Thinking	92.74	86.87	8.67	9.27	8.70	9.00	7.80	8.93	8.90	9.00	8.03	8.57
Claude-Sonnet-4-Thinking	91.03	85.27	8.57	9.00	8.63	8.73	7.57	9.10	8.93	8.70	7.97	8.07
Claude-Sonnet-4	91.00	85.24	8.57	8.80	8.50	8.70	7.80	9.03	8.80	8.80	8.17	8.07
GPT-4o-search	89.40	83.74	8.27	8.77	8.43	8.67	7.77	8.80	8.20	8.73	8.27	7.83
GPT-4o	88.09	82.51	7.90	8.67	8.30	8.33	7.17	8.97	8.57	8.67	7.63	8.30
Gemini-1.5-Pro	85.91	80.47	8.13	8.45	8.30	8.37	7.04	8.17	8.43	8.50	7.48	7.60
o3-mini	84.13	78.80	7.97	8.60	8.30	8.20	6.73	8.57	8.53	7.17	7.03	7.70
Claude-3.5-Sonnet	83.38	78.10	7.97	8.53	8.27	7.93	7.03	8.50	8.00	7.57	6.70	7.60
o1-mini	78.93	73.93	7.27	8.43	7.90	7.53	6.27	8.27	8.17	6.43	6.63	7.03
GPT-4-Turbo	78.57	73.60	6.97	8.17	8.33	7.80	6.00	7.57	8.13	7.00	6.43	7.20
GPT-4-Preview	76.44	71.60	6.90	7.40	8.03	7.30	6.00	7.47	7.63	6.87	6.33	7.67
Baidu-4.0	75.08	70.33	7.27	7.23	7.67	7.43	5.63	6.47	6.80	7.63	7.80	6.40
Baidu-3.5	69.10	64.73	6.20	6.70	7.80	6.83	5.20	5.50	6.00	7.23	6.57	6.70
ChatGLM-Pro	69.10	64.73	5.90	7.07	7.03	7.90	5.43	6.33	5.00	6.67	5.97	7.43
GPT-4-Legacy	66.15	61.96	6.50	6.73	6.60	6.73	5.43	6.10	6.47	5.30	5.20	6.90
Spark-3.0	65.62	61.47	5.77	6.50	7.27	7.30	5.70	5.90	5.03	6.50	5.23	6.27
Claude-3-Haiku	62.93	58.95	5.80	6.60	6.97	6.63	4.83	5.93	6.33	4.80	5.23	5.83
Gemini-Pro	58.18	54.50	4.87	5.43	7.07	6.43	5.10	4.50	4.65	6.33	4.42	5.70
GPT-3.5-turbo	55.42	51.91	4.97	5.37	6.40	6.47	4.43	4.67	5.43	4.20	4.37	5.60
MiniMax-ABAB5	55.33	51.83	3.87	5.63	6.87	6.97	4.33	4.40	2.93	6.13	4.27	6.43

Table 8: Overall and Subject-Level Scores. R_{SOTA}^{model} represents the relative score (0-100 scale) as defined in Equation (2), with Doubao-1.5-Thinking-Pro as the reference SOTA model. S_{model} represents the absolute score (0-100 scale) as defined in Equation (1). Subject-level scores use a 10-point scale.

Model	1000 Questions			Larger Samples		Statistics	
	Trial 1	Trial 2	Trial 3	2000	4000	Mean	Variance
Doubao-1.5-Thinking-Pro	100.0	100.0	100.0	100.0	100.0	100.0	0.00
DeepSeek-V3	96.48	96.87	98.10	98.02	97.53	97.40	0.51
Qwen-3-32B	92.21	92.58	93.93	94.15	93.45	93.26	0.71
GPT-4o	88.08	90.21	91.50	90.69	90.48	90.19	1.63
o3-mini	84.13	85.31	86.69	85.56	86.18	85.57	0.95

Table 9: Ranking stability across different sample sizes. All scores are relative scores with Doubao-1.5-Thinking-Pro as reference (100.0). The three 1000-question trials demonstrate high reproducibility, with low variance indicating robust measurement precision.

Model Name	Round 1			Round 2		
	Doubao	Gemini	GPT-4o	Doubao	Gemini	GPT-4o
<i>Open-source LLMs</i>						
DeepSeek-R1	0.527	0.662	0.960	0.451	0.251	0.977
DeepSeek-V3	0.326	0.627	0.949	0.366	0.147	0.800
Qwen-3-235B	0.486	0.468	0.818	0.232	0.496	0.931
Qwen-3-32B	0.560	0.390	0.886	0.414	0.271	0.872
<i>Closed-source LLMs</i>						
Claude-Sonnet-4	0.309	0.186	0.826	0.402	0.677	0.909
Claude-Sonnet-4-Thinking	0.451	0.399	0.830	0.443	0.373	0.684
Doubao-1.5-Pro	0.629	0.388	0.950	0.383	0.451	0.915
Doubao-1.5-Thinking-Pro	0.707	0.786	0.993	0.745	0.324	0.920
Gemini-2.5-Pro	0.451	0.539	0.831	0.376	0.682	0.858
Gemini-2.5-Pro-Thinking	0.408	0.547	0.843	0.344	0.185	0.859
GPT-4o	0.447	0.507	0.925	0.553	0.417	0.962
o1	0.599	0.535	0.933	0.571	0.555	0.957
o3-mini	0.517	0.397	0.975	0.531	0.381	0.963
Mean	0.493	0.495	0.902	0.447	0.401	0.893

Table 10: Cohen’s κ correlation coefficient between human evaluation and three large language model evaluations across two rounds.