# BEYOND BENCHMARKS: TOWARD CAUSALLY FAITH-FUL EVALUATION OF LARGE LANGUAGE MODELS

## **Anonymous authors**

Paper under double-blind review

#### **ABSTRACT**

Current large language models (LLMs) evaluations overlook that measured LLM performance is produced on a full evaluation system, including many indispensable components, such as workloads, prompting methods, decoding parameters, and the supporting software–hardware stack. Without an explicit, controlled specification of the evaluation system, attributing performance differences to the model itself is unreliable. Our experiments reveal that uncontrolled testing may lead to accuracy variations of up to 70%. To address this urgent issue, we introduce LLM evaluatology, a principled methodology that reduces the evaluation problem to accurately attributing the outcomes to the effect of the evaluated LLM, which is a high-dimensional causal-attribution problem. Empirical results demonstrate that LLM evaluatology not only enhances interpretability and causal validity, but also yields evaluations that are more robust, reproducible, and trustworthy than prevailing benchmarks.

### 1 Introduction

Current LLM evaluation practices are fragmented and ad-hoc, spanning standardized test-style benchmarks (Hendrycks et al., 2020; Huang et al., 2023; Rein et al.; Suzgun et al., 2023; AIME, 2025), human preference-based benchmarks (Chiang et al.; OpenCompass, 2025; Xu et al., 2023), and dynamic or continuously refreshed benchmarks (Jain et al.; Jimenez et al.; White et al.; Zhu et al.; Li et al.). Yet all largely treat the model in isolation, neglecting that measured performance arises from the entire evaluation system, including workloads, prompts, decoding, and even the software-hardware stack. In reality, LLM evaluation is inherently a high-dimensional problem, as these interacting components jointly shape outcomes and complicate attribution. As recent studies show, results can vary sharply with dataset artifacts (Long et al., 2024; Liu et al., 2025), prompt formatting (He et al., 2024), decoding strategies (Shi et al., 2024), or annotator biases (Das et al., 2024). But such analyses remain piecemeal, each targeting a single component without quantifying their combined impact or enabling principled attribution. What is missing is a rigorous methodology that disentangles intrinsic model capability from confounding influences and establishes a reliable foundation for evaluation.

Even under a fully specified evaluation system, LLMs differ fundamentally from traditional single-task or deterministic systems such as conventional algorithms or CPUs. For CPUs, workloads in domains like desktop computing or high-performance computing exhibit well-characterized patterns, allowing evaluation to focus on representative hotspots while treating less common cases as secondary. In contrast, LLM workloads are effectively open-ended: each user can define new tasks across languages, domains, and usage styles. Some tasks resemble those seen during training, others require analogical transformation from familiar patterns, and yet others are entirely novel. This diversity eliminates the notion of a single "typical" workload, making isolated evaluation on a few canonical examples insufficient. In addition, LLMs may produce fluent responses without genuine reasoning or knowledge, so-called hallucinations, meaning that correctly solving one instance does not guarantee mastery of the underlying skill. Consequently, reliable evaluation must consider multiple task variations, from familiar to analogical to novel, in order to disentangle true capability from surface-level correctness. Interpreting performance and attributing capability is therefore both a high-dimensional and a content-sensitive challenge, further amplified by the confounding inherent in the evaluation system.

This paper introduces LLM evaluatology (Fig. 1), a principled methodology for the rigorous evaluation of LLMs based on Evaluatology (Zhan, 2024; Zhan et al., 2024). At its core, we construct a Minimal Evaluation System (MES), which explicitly defines the evaluated object (e.g., standalone LLM or LLM service), the indispensable components influencing performance, and the evaluation conditions (the configuration space formed by admissible settings of these components). By providing a well-defined, controllable system, MES enables systematic exploration of the evaluation configuration space, capturing how different components jointly affect performance and allowing accurate attribution of model capabilities – a solution to the high-dimensional nature of LLM evaluation. To address content sensitivity, we further extend MES into an Augmented MES (A-MES), which transforms existing workloads and generates new instances along semantically related themes. This approach ensures evaluation coverage across three workload layers: workloads the model is likely to have seen, workloads requiring analogical transformation, and entirely novel workloads, thereby mitigating the risks of superficial correctness and hallucination. A-MES offers a structured, reproducible framework that disentangles intrinsic model competence from confounding influences while accommodating the diversity and dynamism of real-world user interactions.

Our experiments reveal several important findings. First, by constructing A-MES, we observe that the accuracy of Doubao varies dramatically with configuration, ranging from 0 to 0.8, highlighting the substantial impact of evaluation settings. Notably, Doubao-1.5-pro ranks first under MES but drops to sixth under A-MES, with a significant gap from the top model, indicating limited generalization ability. Within the Qwen series, we find that the smaller model ranks higher under MES but is surpassed by the larger model under A-MES, suggesting that A-MES provides a more faithful reflection of scaling properties. By contrast, DeepSeek-V3 consistently achieves strong accuracies across all MES and A-MES scenarios, demonstrating the strongest robustness among the tested models. Second, leveraging analysis of variance (ANOVA), xgboost, and linear models, we quantify the impact of each component on model accuracy. All components show measurable influence, with Question Format and COT emerging as the most sensitive, followed by max\_tokens, Shot, and Multi Turn. Furthermore, models exhibit heterogeneous sensitivity to languages: for example, DeepSeek-V3 is most sensitive to Arabic, where its accuracy reaches the lowest among all languages tested. Finally, we validate that our proposed LLM evaluatology provides the closest approximation to the accuracy ground truth, significantly outperforming traditional single-configuration evaluations in reliability and robustness.

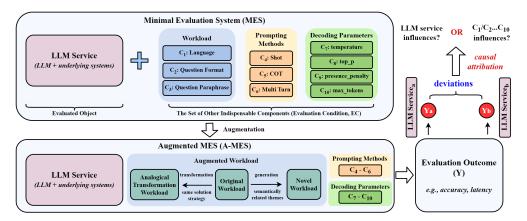


Figure 1: LLM Evaluatology: Measured performance arises from an Augmented Minimal Evaluation System (A-MES), which enables disentangling intrinsic model capability from confounding influences. Here, the evaluation object is defined as the LLM service, comprising the LLM and its underlying systems. When evaluating a standalone LLM, the underlying systems are instead treated as part of the evaluation conditions (EC).

#### 2 Related Work

Broadly, existing benchmarks can be grouped into the following three categories. Standardized test–style benchmarks present problems in the form of test questions, with model outputs compared against reference answers. Representative examples include MMLU (Hendrycks et al., 2020) and its extensions MMLU-Pro (Wang et al., 2024b) and MMLU-Redux (Gema et al., 2025), as well as

C-Eval (Huang et al., 2023) and CMMLU (Li et al., 2024) in the Chinese context. GPQA (Rein et al.) targets graduate-level science, while other datasets focus on specific capabilities such as reasoning (BBH (Suzgun et al., 2023), HellaSwag (Zellers et al., 2019), Winogrande (Sakaguchi et al., 2021)), mathematics (GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al.), AIME (AIME, 2025)), coding (HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), Aider-polyglot (Aider, 2025), MultiPL-E (Cassano et al., 2023)), long-context understanding (L-Eval (An et al., 2024), LongBench (Bai et al., 2024),  $\infty$ Bench (Zhang et al., 2024a), HELMET (Yen et al., 2025)), safety (SafetyBench (Zhang et al., 2024b), Toxigen (Hartvigsen et al., 2022)), instruction-following (IFE-val (Zhou et al., 2023), Multi-Challenge (Sirdeshmukh et al., 2025)), and multimodality (MMBench (Liu et al., 2024), MMMU (Yue et al., 2024), MathVista (Lu et al.)).

Human preference—based benchmarks evaluate models in interactive settings, collecting user judgments instead of relying on fixed test sets. Chatbot Arena (Chiang et al.) is the most prominent example, where pairwise votes are aggregated via Elo ratings. CompassArena (OpenCompass, 2025) and SuperCLUE (Xu et al., 2023) apply similar designs in the Chinese context.

Dynamic or continuously refreshed benchmarks aim to avoid data contamination by relying on newly released or procedurally generated tasks. Examples include LiveCodeBench (Jain et al.) (recent programming contests), SWE-bench (Jimenez et al.) (GitHub issues and PRs), LiveBench (White et al.) (rolling monthly refresh), DyVal (Zhu et al.) (procedural reasoning via DAGs), and Arena-Hard (Li et al.) (real-time crowdsourced challenges).

Table 1: Evaluation Settings on Different Benchmarks (Lang. = Language, Format = Question Format, Para. = Question Paraphrase, M-turn = Multi Turn, Temp. = temperature, PP = presence\_penalty, MaxTok = max\_tokens, ori = original, y = yes, n = no)

Model	Lang.	Format	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
MMLU AIME GPQA MATH SWE-bench IFEval Arena-Hard Human Eval	English English English English English English English English English	ori ori ori ori ori ori ori	n n n n n n n	0/3/5 0 0/5 0/8 0 0 0	y/n y/n y/n y/n y/n y/n y/n y/n	n n n n n n	0.0/0.3/0.5/0.6/0.7 0.0/0.6/0.7 0.4/0.5/0.6/0.7 0.0/0.6/0.7 0.0/0.8 0.0/0.6/0.7 0.0/0.6/0.7 0.0/0.6/0.7	0.8/0.95 0.8/0.95 0.8/0.95 0.8/0.95 0.95 0.8/0.95 0.8/0.95 0.8/0.95	0/1.5 0/1.5 0/1.5 0/1.5 x 0,1.5 x 0,1.5	1024/8192/32768 8192/32768/38912 1024/8192/32768 8192/32768 8192/16384 8192/16384 8192/32768 8192/32768

# 3 MOTIVATION

The flaw of existing LLM evaluation methodology. Existing LLM benchmarks define workload formats and scoring rules, but leave crucial indispensable components uncontrolled, e.g., decoding parameters and prompting methods. As a result, reported evaluation outcomes often do not allow a direct comparison of model differences and may conflate intrinsic capability with arbitrary component settings. To make this issue concrete, we systematically reviewed major benchmarks and compiled a taxonomy of which components are explicitly defined and which are left open (Table 1). Strikingly, many widely used benchmarks, including AIME, specify only a subset of variables while leaving key components underspecified. To quantify the implications, we reconstructed the AIME evaluation space by enumerating plausible settings of uncontrolled components (e.g., COT, temperature, top-p, presence penalty, max tokens), yielding 162 distinct evaluation conditions. Accuracy under these conditions varied by as much as 70% across settings, and the distributions often diverged substantially from the single numbers reported in technical documentation. On some models, the median relative change between our measured accuracy and the accuracy reported in the technical report reached as high as 50%(see Figure 2). Comparable inconsistencies are evident in MMLU (Appendix A.2) and other flagship benchmarks, suggesting that the problem is not dataset-specific but structural across current LLM evaluation methodologies. These findings reveal a fundamental flaw in current practice: a benchmark score is often not a property of the model alone but of the loosely specified evaluation system surrounding it. Without principled control over these confounding components, evaluation becomes unstable, attribution unreliable, and comparisons across models misleading.

The challenges of using Evaluatology for LLM evaluation. Zhan et al. conceptualize evaluation as constructing a minimal system that integrates the evaluation object with indispensable components while considering user requirements (Zhan, 2024; Zhan et al., 2024). Wang et al. illustrate this

approach for CPUs, where a Minimal Evaluation System (MES) isolates CPU behavior from confounding components (Wang et al., 2024a). However, extending Evaluatology to LLMs presents a qualitatively deeper challenge than in the case of CPUs or other deterministic systems. For such conventional artifacts, workloads can be reasonably characterized and stabilized: standardized benchmarks capture dominant usage scenarios and once confounders are controlled, evaluation outcomes largely reflect intrinsic system differences. By contrast, LLM workloads are inherently open-ended and socially constructed, shaped by heterogeneous users, diverse linguistic and cultural contexts, and the continual emergence of novel use cases. In this setting, even the "unit of evaluation" becomes unstable: what qualifies as mainstream, extrapolative, or out-of-distribution shifts across communities and over time. To illustrate, consider the following problem from AIME: "Let A, B, C, and D be points on the hyperbola  $\frac{x^2}{20} - \frac{y^2}{24} = 1$  such that ABCD is a rhombus whose diagonals intersect at the origin. Find the greatest real number that is less than  $BD^2$  for all such rhombi." When evaluated on nine LLMs including deepseek, doubao, gpt series, moonshot, mistral, qwen series, etc.,

five were able to solve this original (seen) workload correctly. However, after performing analogical transformations through inserting distractor: "In a geometric study, we often encounter various shapes and their properties. Also, the concept of symmetry plays an important role in analyzing the relationships between different geometric figures. Let A, B, C, and Dbe points on the hyperbola  $\frac{x^2}{20} - \frac{y^2}{24} = 1$  such that ABCD is a rhombus whose diagonals intersect at the origin. Find the greatest real number that is less than  $BD^2$  for all such rhombi.", none of these models produced correct solutions. This striking contrast illustrates why A-MES is essential: performance on a single workload can be misleading, as models may succeed on problems they have effectively memorized yet fail when the same reasoning must be applied under slightly altered conditions.

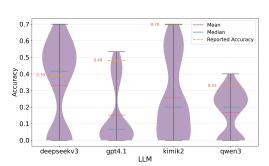


Figure 2: Accuracy deviations on AIME when evaluating with identical workloads across 162 combinations of component settings (COT, temperature, top-p, presence penalty, and max tokens).

#### 4 LLM EVALUATOLOGY

LLM evaluatology consists of three essential steps: (1) defining MES, (2) defining A-MES, and (3) evaluating A-MES and attributing evaluation outcomes.

#### 4.1 DEFINING MINIMAL EVALUATION SYSTEM (MES)

We define the Minimal Evaluation System (MES) for LLM evaluation as the smallest independently runnable system that includes the evaluated object and all indispensable components that materially affect the evaluation outcome. The evaluated object O is not limited to a bare LLM; it can also encompass the broader deployed LLM service that fuses the model with its supporting software and hardware stack. For example, when evaluating through an API, the LLM and its underlying systems should be treated as an inseparable whole, whereas for locally deployed open-source models, the surrounding system environment may either be incorporated into O or explicitly modeled as part of the other indispensable components. Thus, the first step of defining MES is to rigorously define the evaluated object.

The second step in defining MES is to identify the indispensable components that shape evaluation outcomes and to establish their value ranges, collectively denoted as evaluation conditions (EC). We organize EC into three layers, covering workload, prompting method, and decoding parameters, which together yield 10 key factors  $(C_1-C_{10})$ . Workload captures data-related variations, including Language, Question Format, and Question Paraphrase  $(C_1-C_3)$ . Note that Question Paraphrase is introduced as a key component to mitigate hallucination and data contamination, referring to reformulating questions without altering their semantics or correct answers. Prompting method accounts for interaction styles, namely Shot, COT(chain-of-thought), and Multi Turn  $(C_4-C_6)$ . De-

Table 2: Evaluation Conditions: Indispensable Components and Value Ranges

 Language
Question Format
Question Paraphrase
Shot
COT
Multi Turn
temperature

Variable

Chinese, English, Japanese, Arabic, French, Russian Multiple-choice, Fill-in-the-blank Yes, No Yes, No Yes, No 0.0, 1.0, 2.0 0.0, 1.0, 2.0 0.2, 0.6, 1.0 -0.5, 0.5, 1.5 10, 100, 4000

Value Range

coding parameters represent inference controls, including temperature, top\_p, presence\_penalty, and max\_tokens  $(C_7-C_{10})$ . Each component is instantiated with representative values to balance coverage of real-world variability against configuration space tractability. The indispensable components and their value ranges are summarized in Table 2, with both components and their value ranges configurable based on the evaluation object and user-defined requirements. Each MES instance is then specified as  $MES = EC \times O$ , ensuring that performance measurements are attributed correctly while systematically controlling for confounding factors introduced by indispensable components.

## 4.2 Constructing Augmented MES (A-MES)

To further overcome the limitations of traditional evaluation, we extend MES into an augmented form (A-MES) by expanding the workload subspace. Specifically, an MES is defined as  $EC \times O$ , where the evaluation conditions factorize as  $EC = W(workload) \times P(prompting\_methods) \times D(decoding\_parameters)$ . We do augmentation in workload W and leave the non-workload EC components P and D unchanged when building A-MES. Thus  $A-MES = O \times EC_A$ , with  $EC_A = W_A \times P \times D$ ,  $W_A = A(W)$ . A(V) is the augmentation operator that expands the original workload W into an enriched workload W. Practically, A(W) is constructed by partitioning and extending items from the original workload into three purpose-built categories, as shown in Fig. 3:

- Original (Seen) workloads. The workloads within the original benchmark represent cases that the model may have been exposed to during training.
- Analogical transformation (Transformed) workloads. Workloads derived from original
  ones by applying controlled transformations that preserve the underlying solution strategy but may change surface form and answers, including numeric substitutions, distractor
  insertions, and conditional recompositions, where conditional recompositions involve rearranging or swapping problem statements and conditions. These probe the model's ability
  to reason by analogy rather than recall.
- Novel (Out-of-distribution) workloads. Newly created workloads targeting the same concepts and semantically related themes, but unlikely to appear in training corpora. Two complementary strategies are used: (a) recent-source adaptation: harvesting fresh problems (e.g., problems published within a short window unlikely to be in training cutoffs) and adapting them; (b) concept synthesis: generating questions from textbook/academic statements or extracted topic templates to test concept-level mastery.

Why augment only W? LLM evaluation is uniquely sensitive to the space of workloads users present: some queries are seen, some require analogical transfer, others are entirely novel. Expanding the workload subspace in a structured way is therefore the most direct and reproducible means to (a) expose memorization/data-contamination, (b) test analogical/generalization ability, and (c) evaluate true concept transfer — all without conflating these checks with prompt or system confounders.

#### 4.3 EVALUATING ON A-MES AND ATTRIBUTING EVALUATION OUTCOMES

Given the exponentially large space of workload, prompting methods, and decoding parameters, exhaustive testing is generally infeasible. Evaluation on A-MES balances the trade-off between evaluation accuracy and evaluation cost by systematically sampling the configuration space of evaluation conditions. Specifically, for each workload, we perform random sampling of configurations iteratively. For example, after every batch of samples (e.g., 10), we compute the mean performance and confidence interval. Sampling continues until the mean converges within a small threshold (e.g., < 0.002) and the confidence interval length is sufficiently narrow (e.g., < 0.06). Typically, 200–400

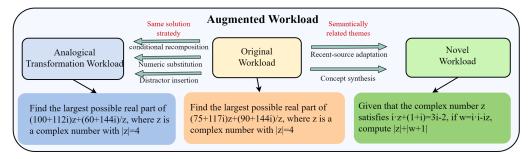


Figure 3: Augment the original workload into analogical transformation and novel workloads.

samples per workload are sufficient. This procedure is justified in two complementary ways: (1) by the law of large numbers, the sample mean converges to the population mean as the number of samples increases; (2) by monitoring both the mean change and confidence interval, we ensure empirical convergence, providing a practical stopping criterion for finite sampling.

The sampled evaluation conditions are then used to test the evaluation object, yielding performance outcomes under diverse settings. One approach to isolate component effects is to use equivalent evaluation conditions, where all component settings are held constant except for the factor of interest; differences in measured performance can thus be attributed directly to that component, effectively mitigating confounding. An alternative and complementary approach is to apply ANOVA (analysis of variance) across the sampled configurations, quantifying the proportion of performance variance explained by each component and enabling systematic attribution of effects. Together, these strategies provide both controlled and statistical means to disentangle intrinsic model capability from the influence of evaluation conditions.

### 5 EVALUATION

In this section, we evaluate the proposed methodology using mainstream LLMs that are publicly accessible, including deepseek-v3, doubao-1.5-pro-32k, gpt-3.5, gpt-4.1, moonshot-v1-8k, mistral-large, mistral-medium, qwen-plus and qwen2.5-32b-instruct. We have three targets. 1) Demonstrate the necessity of constructing MES and A-MES for LLM evaluation by varying the settings of each indispensable component within MES and A-MES. 2) Quantify the contribution of each indispensable component to overall performance variance and identify the key components affecting LLM behavior using ANOVA. 3) Compare LLM evaluatology with traditional LLM evaluation methods and show how it enables accurate attribution of performance differences to specific components.

For online testing, we primarily access the models through their official APIs; however, since the official API for Deepseek v3 has been discontinued, we instead use the API provided by a third-party server deployment. This study employs several widely used and representative benchmark datasets—MMLU, GPQA, and AIME—as the basis for evaluation. Note that due to the page limit, the results of MMLU and GPQA are listed in Appendix A.3. MMLU covers 57 subjects and contains a large collection of multiple-choice questions, widely used to assess models' general knowledge and reasoning abilities. GPQA consists of 448 challenging multiple-choice questions developed and validated by experts in biology, physics, and chemistry, designed to evaluate AI models' reasoning ability on complex scientific problems. AIME is a highly selective U.S. high school mathematics competition, well known for its challenging problems that test deep mathematical reasoning. It is worth noting that our methodology is not tied to any specific benchmark and can be applied to the evaluation of any LLM.

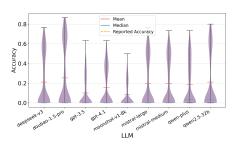
#### 5.1 THE NECESSITY OF CONSTRUCTING MES AND A-MES

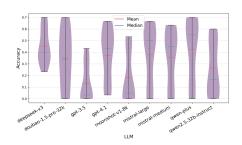
This section demonstrates that LLMs exhibit significant performance variations across different MES and A-MES configurations, thereby underscoring the inadequacy of single-configuration evaluations in accurately capturing their true capabilities.

In the MES experiments, we conducted 500 random samplings without replacement from the MES configuration space described in Section 4.1. The specific components and their corresponding value ranges are summarized in Table 2.

In the A-MES experiments, to verify the effectiveness of the Augmented Minimal Evaluation System (A-MES) proposed in Section 4.2 and comprehensively evaluate the performance of LLMs across diverse task scenarios, we conducted a comparative analysis of their performance on two types of datasets: the original AIME workload and four augmented workloads derived from this original workload. For the analogical transformation workload, we employed two specific methods: the first involves inserting redundant information into the stem of the original question, information that is irrelevant to the problem-solving logic and methods yet consistent with the question scenario, to interfere with the output results of LLMs; The second method involves numeric substitutions. For novel (out-of-distribution) workloads, this study designs two core strategies: the first is a knowledge point-based question generation strategy, which specifically generates new tasks based on the core knowledge points covered in the original questions and combined with the conceptual system and expression paradigm of relevant textbook chapters; The second is an adaptation and transformation strategy based on college entrance examination (gaokao) questions, which involves selecting the latest gaokao questions that match the target knowledge points and generating new tasks by adjusting the scenario of the question, questioning logic, and other aspects.

The experimental results of this study are presented in Figure 4. As shown in Figure 4, significant variations in accuracy trends are observed across different models and configuration spaces. For instance, the accuracy of the deepseek-v3 model fluctuates within a range of 0 to 0.78 under MES experiments, and from 0.23 to 0.7 under A-MES experiments. As shown in Table 3, we have also generated performance rankings for large models based on the original workload, MES workload, and A-MES workload. For the MES and A-MES scenarios, we employ their average accuracy as the performance metric for the LLMs. It is crucial to note that when the accuracy is zero, we sort the models alphabetically based on their names. Drawing insights from the rankings, we observe three key conclusions: first, the original evaluation methodology demonstrates limited effectiveness in benchmarking large language models (LLMs) due to its inability to distinguish performance beyond two models achieving non-zero accuracy scores; second, DeepSeek consistently outperforms all competing models across diverse evaluation conditions, underscoring its robustness and superior generalization capabilities; third, model performance rankings exhibit contextual sensitivity, as evidenced by Doubao's inferior performance relative to DeepSeek in both Original and A-MES workloads, yet its top-ranking achievement in MES, thereby highlighting the non-transitive nature of LLM performance across varying task formulations and data distributions.





(a) Distribution of Model Accuracies on MES

(b) Distribution of Model Accuracies on A-MES

Figure 4: Distribution of Model Accuracies

# 5.2 QUANTIFY THE CONTRIBUTION OF EACH INDISPENSABLE COMPONENT TO OVERALL PERFORMANCE VARIANCE

In LLM evaluation, a key challenge lies in effectively evaluating the contribution of each components illustrated in Fig. 1 to overall performance variance. Given the enormous number of possible EC configuration combinations, exhaustively testing every configuration is computationally infeasible. To address this, we selected a limited number of experimental points from the full space, allowing us to systematically and evenly examine the effects of multiple components and their lev-

378 379

384 385

386 387 388

389 390

391 392 393

396 397

398 399 400

401

407

416 417

426 427

428

429

430

431

Table 3: Performance Rankings of LLMs

Type	1.	2	3	4	5	6	7	8	9
Original	deepseek- v3(0.4)	gpt- 4.1(0.07)	doubao- 1.5-pro- 32k(0)	gpt-3.5(0)	mistral- large(0)	mistral- medium(0)	moonshot- v1-8k(0)	qwen- plus(0)	qwen2.5- 32b- instruct(0)
A-MES	deepseek- v3(0.45)	qwen- plus(0.43)	mistral- large(0.38)	gpt- 4.1(0.37)	mistral- medium(0.36)		qwen2.5- 32b-	moonshot- v1-	gpt- 3.5(0.13)
MES	doubao- 1.5-pro- 32k(0.25)	deepseek- v3(0.21)	qwen2.5- 32b- instruct(0.21)	mistral- large(0.20)	mistral- medium(0.20)	32k(0.34) qwen- plus(0.18)	instruct(0.26) gpt- 4.1(0.16)	8k(0.18) gpt- 3.5(1.10)	moonshot- v1- 8k(0.08)

Note: Models are sorted alphabetically by name when accuracy equals zero.

Table 4: ANOVA results on DeepSeek-V3 (sorted by effect size in descending order)

Factor	Effect Size $\eta^2$	p-value
Question Format	0.399643	0.000
Question Format - COT	0.161394	0.000
ČOT	0.080156	0.000
max_tokens	0.028099	0.000
Question Format - Shot	0.011101	0.000
Language - Question Format	0.008178	0.006
COT - max_tokens	0.006721	0.010
Language - COT	0.004345	0.038
Multi Turn - max_tokens	0.003841	0.050
Language	0.003841	0.046
Shot - max_tokens	0.003669	0.046
Question Format - max_tokens	0.002687	0.100
Language - Multi Turn	0.002600	0.066
temperature - top_p	0.002082	0.178
Question Format - Multi Turn	0.001321	0.244

els on performance with significantly fewer trials. This design reduces experimental cost while maintaining scientific rigor and representativeness.

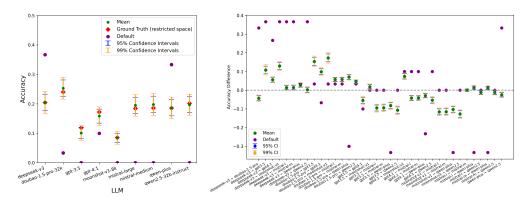
To quantify the proportion of performance variance explained by each MES component, we adopted an analysis of variance (ANOVA) approach. Specifically, for component  $C_1$ – $C_{10}$ , we selected two levels ("high" and "low") within their respective ranges, thereby constructing a subspace of size  $2^{10} = 1024$ . For the Language component, we selected Chinese and English, while for three-valued components we used their maximum and minimum values. Within this subspace, variance decomposition was used to quantify the contributions of different components and their interactions to variations in accuracy. Moreover, we employed a permutation test to evaluate statistical significance, enabling a more robust assessment of component importance without relying on additional distributional assumptions. This procedure yields both the relative importance and the statistical significance of all components.

Taking the DeepSeek-V3 model as an example, Table 4 reports the main effects and two-way interactions that significantly influence its accuracy on the AIME'24 benchmark, with the complete ANOVA results provided in Appendix A.4. Overall, Question type, COT, max\_tokens, and their interactions with other components exhibit the most significant effects. Shot, Multi turn, and Language also show significant effects, while the remaining components have only limited impact.

Consistent patterns were observed across other LLMs (see Appendix A.4). Using p < 0.05 as the significance threshold, we found that the main effects of Question format and COT, or their interactions with other components, were consistently significant across all LLMs. Furthermore, max\_tokens, Shot, and Multi turn also reached significance for the vast majority of models. In addition to these five core components, Language, top\_p, and temperature were significant for some models. It is worth noting that for the remaining two components, Question Paraphrase and presence\_penalty, the p-values did not meet the significance threshold, but reached 0.19 and 0.16, respectively, on GPT-4.1. This suggests that they may exert some influence on model performance, although the evidence is not sufficient for a definitive conclusion.

# COMPARE LLM EVALUATOLOGY WITH TRADITIONAL LLM EVALUATION METHODS AND ATTRIBUTE THE PERFORMANCE DIFFERENCES TO SPECIFIC COMPONENTS

This section demonstrates that evaluating models under a single configuration fails to capture their true capabilities, while LLM evaluatology not only yields results in strong agreement with the ground truth, but also attributes the performance differences to specific components.



(a) Accuracy confidence intervals (b) Accuracy difference confidence intervals of different models on of different LLMs on AIME

AIME

Figure 5: Comparison of LLM Evaluatology and Traditional Method on AIME

Based on the randomly sampled data collected from the complete configuration space, we estimated the overall average accuracies of different models on the same benchmark using their 95% and 99% confidence intervals. As illustrated in Figure 5a, we report the performance of different models on AIME'2024, where the purple dots denote the test results under the commonly adopted default setting, using the original workloads without optimized prompting methods and with default decoding parameters. It can be observed that the purple dots are far from the confidence intervals (interval estimation of the population mean) obtained through random sampling, showing that evaluating a model under a single configuration is unreliable.

Figure 5 further presents, on AIME'2024, the accuracy differences between two models under the default configuration, together with the 95% and 99% confidence intervals constructed from accuracy differences observed across sampled equivalent evaluation configurations. In 9 cases, the confidence intervals and the default accuracy differences fall on opposite sides of the zero line, revealing contradictions in the conclusions regarding model superiority. For instance, the 99% confidence interval for the mean accuracy difference between Doubao-1.5-pro and GPT-4.1 lies entirely above the zero line, implying that overall Doubao-1.5-pro outperforms GPT-4.1. However, if one were to rely on the result of a single experiment under the default configuration, the accuracy difference would fall below the zero line, leading instead to the opposite conclusion that GPT-4.1 outperforms Doubao-1.5-pro. This "conclusion reversal" highlights the limitations of relying solely on single-configuration testing. More detailed results on additional benchmarks including MMLU and GPQA can be found in the Appendix.

Furthermore, we selected the five most influential components for a cost-efficient accuracy test on each LLM, based on the ANOVA data in Section 5.2. We then constructed the configuration subspace restricted to these components and conducted exhaustive testing within this subspace. The mean performance obtained was taken as a "restricted-space ground truth." As shown by the red diamond in Figure 5a, for all models, this reference truth fell within the confidence intervals estimated from random sampling, thereby demonstrating both the validity and the robustness of the proposed LLM evaluatology method.

# 6 Conclusion

LLM Evaluationsy establishes a principled methodology for assessing LLMs through an Augmented Minimal Evaluation System (A-MES), explicitly accounting for both intrinsic model capabilities and the many confounding components that shape observed performance, thereby enabling accurate attribution of performance differences to their true sources. Our analysis reveals that meaningful evaluation of LLMs requires careful consideration of both workload heterogeneity and the vast space of evaluation condition (EC) configurations. We advocate for the adoption of evaluatology as a foundational paradigm, encouraging the community to develop richer workload augmentation strategies and robust evaluation practices that mirror the complexity of actual deployment scenarios.

# REFERENCES

- Aider. Aider Ilm leaderboards. https://aider.chat/docs/leaderboards/, 2025. Accessed: 2025-09-13.
- AIME. AIME Problems and Solutions. https://artofproblemsolving.com/wiki/index.php/AIME\_Problems\_and\_Solutions, 2025. Accessed: 2025-09-05.
  - Chenxin An, Shansan Gong, Ming Zhong, Xingjian Zhao, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. L-eval: Instituting standardized evaluation for long context language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pp. 14388–14411, 2024.
  - Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
  - Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3119–3137, 2024.
  - Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. Multiple: A scalable and polyglot approach to benchmarking neural code generation. *IEEE Transactions on Software Engineering*, 49(7):3675–3691, 2023.
  - Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
  - Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. In *Forty-first International Conference on Machine Learning*.
  - 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.
  - Amit Das, Zheng Zhang, Najib Hasan, Souvika Sarkar, Fatemeh Jamshidi, Tathagata Bhattacharya, Mostafa Rahgouy, Nilanjana Raychawdhary, Dongji Feng, Vinija Jain, et al. Investigating annotator bias in large language models for hate speech detection. *arXiv preprint arXiv:2406.11109*, 2024.
  - Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, et al. Are we done with mmlu? In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 5069–5096, 2025.
  - Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pp. 3309–3326, 2022.
- Jia He, Mukund Rungta, David Koleczek, Arshdeep Sekhon, Franklin X Wang, and Sadid Hasan.

  Does prompt formatting have any impact on llm performance? *arXiv preprint arXiv:2411.10541*, 2024.
- 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)*.

- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2020.
  - Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36: 62991–63010, 2023.
  - Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. In *The Thirteenth International Conference on Learning Representations*.
  - Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*.
  - Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. Cmmlu: Measuring massive multitask language understanding in chinese. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 11260–11285, 2024.
  - Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. In *Forty-second International Conference on Machine Learning*.
  - Jiacheng Liu, Mayi Xu, Qiankun Pi, Wenli Li, Ming Zhong, Yuanyuan Zhu, Mengchi Liu, and Tieyun Qian. Format as a prior: Quantifying and analyzing bias in llms for heterogeneous data. *arXiv* preprint arXiv:2508.15793, 2025.
  - Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pp. 216–233. Springer, 2024.
  - Do Xuan Long, Hai Nguyen Ngoc, Tiviatis Sim, Hieu Dao, Shafiq Joty, Kenji Kawaguchi, Nancy F Chen, and Min-Yen Kan. Llms are biased towards output formats! systematically evaluating and mitigating output format bias of llms. *arXiv* preprint arXiv:2408.08656, 2024.
  - Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*.
  - OpenCompass. Compassarena. https://opencompass.org.cn/arena, 2025. Accessed: 2025-09-08.
  - David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
  - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
  - Chufan Shi, Haoran Yang, Deng Cai, Zhisong Zhang, Yifan Wang, Yujiu Yang, and Wai Lam. A thorough examination of decoding methods in the era of llms. *arXiv preprint arXiv:2402.06925*, 2024.
  - Ved Sirdeshmukh, Kaustubh Deshpande, Johannes Mols, Lifeng Jin, Ed-Yeremai Cardona, Dean Lee, Jeremy Kritz, Willow Primack, Summer Yue, and Chen Xing. Multichallenge: A realistic multi-turn conversation evaluation benchmark challenging to frontier llms. *arXiv preprint* arXiv:2501.17399, 2025.

- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. In *ACL* (*Findings*), 2023.
- Chenxi Wang, Lei Wang, Wanling Gao, Yikang Yang, Yutong Zhou, and Jianfeng Zhan. Achieving consistent and comparable cpu evaluation outcomes. *arXiv* preprint arXiv:2411.08494, 2024a.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging multitask language understanding benchmark. *Advances in Neural Information Processing Systems*, 37:95266–95290, 2024b.
- Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Benjamin Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, et al. Livebench: A challenging, contamination-limited llm benchmark. In *The Thirteenth International Conference on Learning Representations*.
- Liang Xu, Anqi Li, Lei Zhu, Hang Xue, Changtai Zhu, Kangkang Zhao, Haonan He, Xuanwei Zhang, Qiyue Kang, and Zhenzhong Lan. Superclue: A comprehensive chinese large language model benchmark. *arXiv preprint arXiv:2307.15020*, 2023.
- Zhengxin Yang, Wanling Gao, Chunjie Luo, Lei Wang, and Jianfeng Zhan. Quality at the tail. CoRR, abs/2212.13925, 2022. doi: 10.48550/ARXIV.2212.13925. URL https://doi.org/10.48550/arXiv.2212.13925.
- Howard Yen, Tianyu Gao, Minmin Hou, Ke Ding, Daniel Fleischer, Peter Izsak, Moshe Wasserblat, and Danqi Chen. Helmet: How to evaluate long-context language models effectively and thoroughly. International Conference on Learning Representations (ICLR), 2025.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multi-modal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, 2019.
- Jianfeng Zhan. A short summary of evaluatology: The science and engineering of evaluation, 2024.
- Jianfeng Zhan, Lei Wang, Wanling Gao, Hongxiao Li, Chenxi Wang, Yunyou Huang, Yatao Li, Zhengxin Yang, Guoxin Kang, Chunjie Luo, et al. Evaluatology: The science and engineering of evaluation. *BenchCouncil Transactions on Benchmarks, Standards and Evaluations*, 4(1):100162, 2024.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, et al. ∞bench: Extending long context evaluation beyond 100k tokens. *arXiv preprint arXiv:2402.13718*, 2024a.
- Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. Safetybench: Evaluating the safety of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pp. 15537–15553, 2024b.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv* preprint arXiv:2311.07911, 2023.
- Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval: Dynamic evaluation of large language models for reasoning tasks. In *The Twelfth International Conference on Learning Representations*.

# A APPENDIX

#### A.1 THE USE OF LARGE LANGUAGE MODELS (LLMS)

During the manuscript preparation, we leveraged large language models (LLMs) to assist in refining and polishing the text. Specifically, the LLM was used to improve sentence clarity and enhance linguistic fluency, while all scientific content, reasoning, and results were independently authored and verified by the researchers. This approach facilitated more concise and readable presentation without affecting the technical accuracy.

### A.2 EVALUATION SETTING ON DIFFERENT BENCHMARKS

Table 5: Evaluation Settings Reported in Technical Reports of Different LLMs. (Lang. = Language, Q-type = Question type, Para. = Paraphrase, M-turn = Multi turn, Temp. = Temperature, PP = presence\_penalty, MaxTok = max\_tokens)

(a	) Eva	aluation	Settings	on	AIME	3'2024
----	-------	----------	----------	----	------	--------

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
DeepSeek-R1	english	origin	origin	0	yes	х	0.6	0.95	х	32768
DeepSeek-V3	english	origin	origin	X	X	X	0.7	X	X	8192
Kimî K2	english	origin	origin	X	no	X	0.0	fixed	X	8192
Kimi K1.5	english	origin	origin	X	yes	x	X	X	X	X
Owen2	-			No	ot evalua	ted on AIM	E			
Öwen2.5				No	ot evalua	ted on AIM	E			
Qwen3	english	origin	origin	X	no	x	0.7	0.8	1.5	32768
GPT-4	-			No	ot evalua	ted on AIM	E			
GPT-4.1	english	origin	origin	X	X	x	X	X	X	X
GPT-5	-			No	ot evalua	ted on AIM	E			
Claude Opus 4				No	ot evalua	ted on AIM	E			
Mistral Small3.1				No	ot evalua	ted on AIM	E			
Mistral Medium3				No	ot evalua	ted on AIM	E			
Mistral Large2				No	ot evalua	ted on AIM	E			

#### (b) Evaluation Settings on MMLU

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
DeepSeek-R1	english	origin	origin	0	yes	х	0.5	х	х	1024
DeepSeek-V3	english	origin	origin	0	yes	X	0.5	X	X	1024
Kimi K2	english	origin	origin	X	no	x	0.0	fixed	X	8192
Kimi K1.5	english	origin	origin	X	yes	X	x	X	X	X
Qwen2	english	origin	origin	5	x	X	x	X	X	X
Owen2.5	english	origin	origin	5	X	X	x	X	X	X
Qwen3	english	origin	origin	5	X	X	x	X	X	X
GPT-4	multiple	origin	origin	5/3	no	X	X	X	X	X
GPT-4.1	multiple	origin	origin	X	X	X	x	X	X	X
GPT-5	•		-	N	ot evaluated	on MMLU				
Claude Opus 4	multiple	origin	origin	X	yes/no	X	X	X	X	X
Mistral Small3.1	english	origin	origin	X	x	X	x	X	X	X
Mistral Medium3	Č		-	N	ot evaluated	on MMLU				
Mistral Large2	multiple	origin	origin	X	X	X	X	X	X	X

#### (c) Evaluation Settings on GPQA

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
DeepSeek-R1	english	origin	origin	0	yes	х	0.5	х	х	1024
DeepSeek-V3	english	origin	origin	0	yes	x	0.5	X	x	1024
Kimî K2	english	origin	origin	X	no	X	0.0	fixed	x	8192
Kimi K1.5	english	origin	origin	X	yes	x	x	X	X	x
Owen2	english	origin	origin	X	x	X	X	X	x	X
Owen2.5	english	origin	origin	X	X	X	X	X	x	X
Öwen3	english	origin	origin	X	yes/no	X	0.6/0.7	0.95/0.8	0/1.5	32768
GPT-4	Č	C	C		Not evalu	ated on GPO	)A			
GPT-4.1	english	origin	origin	X	X	X	X	X	x	X
GPT-5	english	origin	origin	X	0/1	X	X	X	x	X
Claude Opus 4	english	origin	origin	X	0/1	x	x	x	x	X
Mistral Small3.1	english	origin	origin	X	X	x	x	x	x	x
Mistral Medium3	english	origin	origin	5	1	x	x	x	x	x
Mistral Large2	3				Not evalu	ated on GPC	QA			

# (d) Evaluation Settings on MATH

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
DeepSeek-R1	english	origin	origin	0/8	yes/no	х	0	х	х	32768
DeepSeek-V3	english	origin	origin	0/8	yes/no	X	0	X	X	8192
Kimi K2	english	origin	origin	X	no	X	0.0	fixed	X	8192
Kimi K1.5	english	origin	origin	X	yes	x	x	X	x	x
Owen2	english	origin	origin	X	x	X	X	X	X	X
Qwen2.5	english	origin origin	origin	X	X	X	x	X	X	X
Qwen3 GPT-4 GPT-4.1 GPT-5 Claude Opus 4	english	origin	origin	х	Not evalua Not evalua	x ated on MA' ated on MA' ated on MA' ated on MA'	TH TH	0.95/0.8	0/1.5	32768
Mistral Small3.1	english	origin	origin	X	X	X	X	X	X	X
Mistral Medium3	english	origin	origin	0	0	X	x	X	X	x
Mistral Large2	english	origin	origin	0	0	Х	X	Х	X	х

# 

# 

# 

# 

# 

# 

# 

# 

# 

# 

# (e) Evaluation Settings on SWE-bench

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
DeepSeek-R1	english	origin	origin	х	х	х	0.8	х	х	x
DeepSeek-V3	english	origin	origin	X	X	X	0.8	X	X	X
Kimi K2	english	origin	origin	x	no	X	0.0	fixed	X	8192/16384
Kimi K1.5 Qwen2 Qwen2.5(pre) Qwen3(pre) GPT-4				No No No	t evaluate t evaluate t evaluate	ed on SWE-bed	ench ench ench			
GPT-4.1	english	origin	origin	X	X	X	X	X	X	X
GPT-5	english	origin	origin	X	0/1	X	X	X	X	X
Claude Opus 4	english	origin	origin	X	0/1	X	X	0.95	X	X
Mistral Small3.1 Mistral Medium3 Mistral Large2				No	t evaluate	ed on SWE-b ed on SWE-b ed on SWE-b	ench			

# (f) Evaluation Settings on IFEval

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTo
DeepSeek-R1	english	origin	origin	0	0	0	х	х	х	
DeepSeek-V3	english	origin	origin	0	0	0	X	X	X	X
Kimi K2	english	origin	origin	X	no	x	0.0	fixed	x	8192
Kimi K1.5	english	origin	origin	X	yes	x	X	X	X	X
Qwen2	english	origin	origin	X	x	X	X	X	X	X
Owen2.5	english	origin	origin	X	X	X	X	X	X	X
Qwen3 GPT-4	english	origin origin	origin	x	yes/no Not evalu	x ated on IFE	0.6/0.7 val	0.95/0.8	0/1.5	32768
GPT-4.1 GPT-5 Claude Opus 4	english	origin	origin	х	Not evalu	x ated on IFE ated on IFE	/al	х	х	x
Mistral Small3.1 Mistral Medium3 Mistral Large2	english	origin	origin	0	0	ated on IFE x ated on IFE	x	x	x	x

# (g) Evaluation Settings on Arena-Hard

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTol
DeepSeek-R1	english	origin	origin	0	0	0	config	default	default	user-set
DeepSeek-V3	english	origin	origin	0	0	0	config	default	default	user-set
Kimi K2	english	origin	origin	X	no	x	0.0	fixed	x	8192
Kimi K1.5	Č				Not evaluat	ed on Arena	-Hard			
Qwen2	english	origin	origin	X	X	X	X	X	x	x
Qwen2.5	english	origin	origin	X	X	X	X	X	X	X
Qwen3	english	origin	origin	X	yes/no	X	0.6/0.7	0.95/0.8	0/1.5	32768
GPT-4	-	_	_		Not evaluat	ted on Arena	-Hard			
GPT-4.1						ted on Arena				
GPT-5						ted on Arena				
Claude Opus 4						ted on Arena				
Mistral Small3.1						ted on Arena				
Mistral Medium3					Not evaluat	ted on Arena	-Hard			
Mistral Large2	english	origin	origin	X	X	X	X	X	X	X

# (h) Evaluation Settings on HumanEval

Model	Lang.	Q-type	Para.	Shot	COT	M-turn	Temp.	top_p	PP	MaxTok
DeepSeek-R1	english	origin	origin	0	0	0	varied	0.95	х	32768
DeepSeek-V3	english	origin	origin	0	0	0	varied	0.95	X	8192
Kimi K2	-	-	-	Not	evaluated	on HumanE	val			
Kimi K1.5	english	origin	origin	X	yes	X	X	X	X	X
Qwen2	english	origin	origin	X	X	X	x	X	X	x
Qwen2.5	english	origin	origin	X	X	X	X	X	X	X
Qwen3	-		-	Not	evaluated	on HumanE	val			
GPT-4	english	origin	origin	0	0	X	0.3	X	X	X
GPT-4.1	-	-	-			on HumanE				
GPT-5				Not	evaluated	on HumanE	val			
Claude Opus 4				Not	evaluated	on HumanE	val			
Mistral Small3.1	english	origin	origin	X	X	X	X	X	X	X
Mistral Medium3	english	origin	origin	0	0	X	X	X	X	X
Mistral Large2	english	origin	origin	x	X	X	X	X	Х	х

# A.3 THE RESULT ON MMLU, GPQA

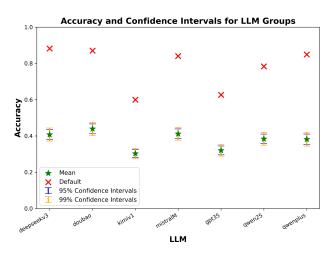
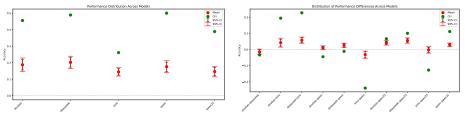


Figure 6: Accuracy confidence intervals of different LLMs on MMLU



(a) Accuracy confidence intervals (b) Accuracy difference confidence intervals of different models on GPQA GPQA

Figure 7: Comparison of LLM Evaluatology and Traditional Testing on GPQA

# A.4 ANOVA ANALYSIS RESULTS ON DIFFERENT LLMS

Table 6: Complete ANOVA results on LLMs (sorted by effect size in descending order)

# (a) Complete ANOVA results on DeepSeek-V3

# (b) Complete ANOVA results on Doubao-1.5-pro-32k

Factor	Effect Size $\eta^2$	p-value	Factor	Effect Size $\eta^2$	p-value
Question Format	0.399643	0.000	Question Format	0.467626	0.000
Question Format-COT	0.161394	0.000	Question Format-COT	0.259657	0.000
COT	0.080156	0.000	COT	0.130549	0.000
max_tokens	0.028099	0.000	max_tokens	0.009192	0.002
Question Format-Shot	0.011101	0.000	COT-max_tokens	0.008943	0.006
Language-Question Format	0.008178	0.006	max_tokens-presence_penalty	0.000671	0.388
COT-max_tokens	0.006721	0.010	Shot	0.000638	0.408
Language-COT	0.004345	0.038	Shot-Multi Turn	0.000605	0.424
Multi Turn-max_tokens	0.003841	0.050	Question Paraphrase-max_tokens	0.000502	0.466
Language	0.003841	0.046	Multi Turn	0.000473	0.466
Shot-max_tokens	0.003669	0.046	Multi Turn-max_tokens	0.000464	0.472
Ouestion Format-max_tokens	0.002687	0.100	Language-temperature	0.000455	0.468
Language-Multi Turn	0.002600	0.066	COT-Multi Turn	0.000427	0.496
temperature-top_p	0.002000	0.178	Question Format-Shot	0.000427	0.546
Question Format-Multi Turn	0.002032	0.244	Question Format-max_tokens	0.000400	0.538
Question Paraphrase	0.001321	0.252		0.000392	0.534
	0.001240	0.252	temperature	0.000392	0.530
Question Format-temperature			Question Format-temperature		
Shot-temperature	0.000926	0.326	Question Paraphrase-presence_penalty	0.000349	0.552
Multi Turn	0.000793	0.364	Language-top_p	0.000326	0.584
temperature	0.000587	0.396	temperature-top_p	0.000326	0.588
COT-Multi Turn	0.000587	0.466	Language-presence_penalty	0.000310	0.578
COT-top_p	0.000573	0.482	Language-Multi Turn	0.000295	0.562
Language-Shot	0.000534	0.466	Multi Turn-presence_penalty	0.000272	0.604
presence_penalty	0.000435	0.488	Language-COT	0.000265	0.586
Question Format-Question Paraphrase	0.000411	0.562	Shot-COT	0.000224	0.628
Question Paraphrase-max_tokens	0.000207	0.626	Question Paraphrase-top_p	0.000205	0.642
Language-Question Paraphrase	0.000198	0.660	Question Format-Multi Turn	0.000158	0.704
Shot-COT	0.000161	0.646	COT-presence_penalty	0.000147	0.696
COT-temperature	0.000140	0.698	max_tokens-top_p	0.000147	0.674
Language-max_tokens	0.000140	0.712	presence_penalty	0.000117	0.730
COT-presence_penalty	0.000140	0.668	temperature-max_tokens	0.000108	0.740
Shot	0.000134	0.706	Question Format-presence_penalty	0.000090	0.764
Question Format-presence_penalty	0.000133	0.708	Question Format-Question Paraphrase	0.000067	0.792
Shot-top_p	0.000109	0.736	Question Paraphrase	0.000054	0.836
max_tokens-presence_penalty	0.000092	0.768	Shot-top_p	0.000047	0.818
max_tokens-top_p	0.000086	0.790	Shot-presence_penalty	0.000047	0.792
top_p	0.000058	0.812	Language	0.000047	0.814
Shot-Multi Turn	0.000054	0.802	Language-Question Paraphrase	0.000043	0.840
temperature-max_tokens	0.000034	0.836	Multi Turn-temperature	0.000036	0.832
Language-presence_penalty	0.000038	0.854	COT-top_p	0.000036	0.854
	0.000033	0.902		0.000034	0.834
top_p-presence_penalty		0.874	COT-temperature		
Multi Turn-temperature	0.000029		Language-Shot	0.000034	0.850
Multi Turn-top_p	0.000026	0.884	Shot-max_tokens	0.000024	0.864
Question Paraphrase-top_p	0.000018	0.898	temperature-presence_penalty	0.000015	0.906
temperature-presence_penalty	0.000013	0.902	top_p-presence_penalty	0.000013	0.902
Question Paraphrase-presence_penalty	0.000011	0.910	Question Paraphrase-COT	0.000007	0.908
Language-top_p	0.000008	0.942	Language-Question Format	0.000002	0.956
Question Paraphrase-Shot	0.000006	0.930	Question Paraphrase-temperature	0.000001	0.960
Multi Turn-presence_penalty	0.000004	0.946	Shot-temperature	0.000001	0.962
Language-temperature	0.000004	0.940	Question Paraphrase-Multi Turn	0.000001	0.966
Question Format-top_p	0.000003	0.962	Language-max_tokens	0.000001	0.970
Question Paraphrase-Multi Turn	0.000003	0.972	Question Paraphrase-Shot	0.000001	0.980
Shot-presence_penalty	0.000001	0.966	top_p	0.000001	0.970
Question Paraphrase-COT	0.000001	0.996	Question Format-top_p	0.000000	0.984
Question Paraphrase-temperature	0.000000	0.994	Multi Turn-top_p	0.000000	0.992

# (c) Complete ANOVA results on GPT-3.5

### (d) Complete ANOVA results on GPT-4.1

Factor	Effect Size $(\eta^2)$	p-value
Question Format	0.417428	0.00000
Question Format-COT	0.199209	0.00000
COT	0.199208	0.00000
temperature	0.006046	0.01600
Question Format-temperature	0.005781	0.02000
Shot-Multi Turn	0.005715	0.01000
Shot-COT	0.005651	0.02600
max_tokens	0.005396	0.02400
Question Format-max_tokens	0.004434	0.04400
temperature-top_p	0.004320	0.05200
top_p	0.003119	0.07400
Question Format-top_p	0.002570	0.10400
COT-max_tokens	0.001773	0.17000
Question Format-Shot	0.000853	0.34400
COT-Multi Turn	0.000828	0.35200
Language-Multi Turn	0.000779	0.36400
Language Mana Tam	0.000687	0.42800
Language-Question Format	0.000600	0.48400
COT-top_p	0.000443	0.49600
		0.50800
COT-temperature	0.000425	
COT-presence_penalty	0.000408	0.50800
Question Paraphrase-Multi Turn	0.000407	0.53400
Question Paraphrase-COT	0.000391	0.50600
Shot	0.000296	0.57800
Shot-temperature	0.000281	0.59800
Language-Shot	0.000253	0.59800
max_tokens-presence_penalty	0.000201	0.62800
Question Paraphrase-top_p	0.000177	0.69600
Shot-max_tokens	0.000135	0.72000
Language-top_p	0.000115	0.69000
Question Paraphrase-presence_penalty	0.000115	0.71400
Language-temperature	0.000106	0.75000
Language-max_tokens	0.000106	0.74000
Shot-presence_penalty	0.000089	0.78600
Multi Turn-temperature	0.000081	0.78400
Language-COT	0.000074	0.76600
Multi Turn-presence_penalty	0.000074	0.78800
Question Format-Question Paraphrase	0.000067	0.77600
Question Paraphrase	0.000053	0.80800
Question Format-Multi Turn	0.000047	0.82600
presence_penalty	0.000047	0.85400
temperature-max_tokens	0.000036	0.85800
Question Paraphrase-max_tokens	0.000036	0.85400
temperature-presence_penalty	0.000031	0.88400
Language-Question Paraphrase	0.000027	0.85800
Question Format-presence_penalty	0.000027	0.90000
Question Paraphrase-temperature	0.000027	0.88200
Multi Turn	0.000027	0.87200
Multi Turn-max_tokens	0.000018	0.91600
Question Paraphrase-Shot	0.000009	0.92000
	0.000009	0.92000
Multi Turn-top_p		0.94200
	0.000007	0.54600
	0.000007	0.04204
top_p-presence_penalty Shot-top_p Language-presence_penalty	0.000007 0.000002	0.94200

(u) Complete ANOVA I	zsuits on Oi 1-	<del></del>
Factor	Effect Size $(\eta^2)$	p-value
Question Format	0.289086	0.000000
Question Format-COT	0.180162	0.000000
COT-max_tokens	0.054845	0.000000
max_tokens	0.053399	0.000000
COT	0.027181	0.000000
temperature-top_p	0.006685	0.010000
Question Format-Shot	0.006529	0.030000
temperature	0.005619	0.016000
Shot	0.004963	0.028000
max_tokens-top_p	0.004019	0.044000
Language-max_tokens	0.003907	0.062000
Question Format-temperature	0.003732	0.072000
top_p	0.003364	0.098000
temperature-max_tokens	0.002990	0.140000
Question Paraphrase-Shot	0.002362	0.160000
Question Paraphrase-presence_penalty	0.001939	0.194000
Shot-Multi Turn	0.001842	0.194000
COT-temperature	0.001708	0.230000
Shot-COT	0.001589	0.290000
Language-COT	0.001517	0.246000
COT-presence_penalty	0.001406	0.276000
Question Format-max_tokens	0.001360	0.306000
temperature-presence_penalty	0.001284	0.292000
Language-Shot	0.001276	0.262000
Language-Question Paraphrase	0.001270	0.324000
presence_penalty	0.001184	0.362000
COT-top_p	0.001183	0.300000
COT-Multi Turn	0.001103	0.376000
Multi Turn-max_tokens	0.000988	0.376000
Shot-top_p	0.000754	0.422000
Question Paraphrase-max_tokens	0.000703	0.438000
Multi Turn-top_p	0.000367	0.574000
Question Format-presence_penalty	0.000348	0.584000
Question Format-Question Paraphrase	0.000348	0.584000
Language-presence_penalty	0.000339	0.586000
max_tokens-presence_penalty	0.000324	0.600000
top_p-presence_penalty	0.000324	0.618000
Shot-presence_penalty	0.000192	0.702000
Shot-max_tokens	0.000164	0.724000
Question Paraphrase	0.000146	0.746000
Question Format-top_p	0.000140	0.728000
Shot-temperature	0.000134	0.724000
Question Format-Multi Turn	0.000133	0.718000
Question Paraphrase-top_p	0.000132	0.798000
Multi Turn-presence_penalty	0.000055	0.828000
Language-top_p	0.000033	0.896000
Language-Question Format	0.000023	0.910000
	0.000022	0.886000
Multi Turn-temperature	0.000020	0.922000
Question Paraphrase-temperature Language-Multi Turn	0.000010	0.922000
	0.000009	0.922000
Question Paraphrase-COT	0.000009	0.922000
Question Paraphrase-Multi Turn		0.918000
Language Multi Turn	0.000001	
Multi Turn	0.000000	0.996000 0.998000
Language-temperature	0.000000	0.998000

### (e) Complete ANOVA results on Qwen2.5

#### (f) Complete ANOVA results on Qwen Plus

Factor	Effect Size $(\eta^2)$	p-valu
Question Format	0.454352	0.00000
Question Format-COT	0.204235	0.00000
COT	0.200224	0.0000
Question Format-Shot	0.009265	0.00000
Shot	0.007983	0.00800
Shot-COT	0.006715	0.00200
Multi Turn	0.003717	0.04600
Question Format-Multi Turn	0.003657	0.05200
max_tokens	0.002764	0.08000
Language-Question Format	0.002764	0.1280
COT-max_tokens	0.001865	0.16400
Language	0.001720	0.1960
COT-Multi Turn	0.001720	0.2140
Multi Turn-max_tokens	0.001110	0.32200
	0.0001110	0.3220
Question Format-max_tokens	0.000922	0.29200
Language-COT		0.3820
Question Paraphrase-presence_penalty	0.000710	
Language-Multi Turn	0.000538	0.4780
temperature	0.000430	0.44400
Shot-temperature	0.000380	0.53400
Question Format-temperature	0.000371	0.55000
Language-Question Paraphrase	0.000343	0.55000
temperature-presence_penalty	0.000334	0.56600
Language-temperature	0.000290	0.60200
top_p	0.000257	0.60000
max_tokens-top_p	0.000242	0.6180
Shot-Multi Turn	0.000234	0.65400
COT-temperature	0.000219	0.66000
Question Format-Question Paraphrase	0.000165	0.6740
Question Paraphrase-top_p	0.000165	0.6720
Language-max_tokens	0.000165	0.67400
Question Paraphrase	0.000118	0.7420
Language-top_p	0.000107	0.76000
temperature-top_p	0.000107	0.72400
Multi Turn-presence_penalty	0.000083	0.7620
max_tokens-presence_penalty	0.000075	0.78400
Shot-max_tokens	0.000062	0.77000
presence_penalty	0.000051	0.83000
Question Format-top_p	0.000038	0.82400
top_p-presence_penalty	0.000038	0.83000
top_p-presence_penalty Multi Turn-temperature	0.000029	0.85600
Question Paraphrase-temperature	0.000027	0.8820
COT-presence_penalty	0.000022	0.84600
Multi Turn-top_p	0.000022	0.8820
Question Format-presence_penalty	0.000022	0.86600
Language-presence_penalty	0.000012	0.91200
temperature-max_tokens	0.000012	0.9180
Question Paraphrase-Multi Turn	0.000009	0.94200
Language-Shot	0.000007	0.92600
Question Paraphrase-Shot	0.000005	0.92800
COT-top_p	0.000003	0.95400
Shot-presence_penalty	0.000002	0.96000
Shot-top_p	0.000001	0.98200
Question Paraphrase-COT	0.000000	0.9840
Question Paraphrase-max <sub>t</sub> okens	0.000000	0.99000

(1) Complete ANOVA lesuits on Qwell Flus			
Factor	Effect Size $(\eta^2)$	p-value	
Question Format	0.302717	0.000000	
Question Format-COT	0.259678	0.000000	
COT	0.098747	0.000000	
Shot	0.047447	0.000000	
Question Format-Shot	0.042448	0.000000	
Shot-COT	0.024956	0.000000	
COT-max_tokens	0.016596	0.000000	
max_tokens	0.016280	0.000000	
Language	0.005019	0.024000	
Language-Question Format	0.003272	0.060000	
temperature-top_p	0.002324	0.112000	
Multi Turn-max_tokens	0.001979	0.154000	
Question Format-temperature	0.001662	0.202000	
top_p	0.001282	0.250000	
Language-Shot	0.000991	0.296000	
Question Paraphrase-Multi Turn	0.000877	0.320000	
COT-temperature	0.000822	0.368000	
Multi Turn	0.000736	0.428000	
Question Format-Question Paraphrase	0.000654	0.458000	
temperature	0.000608	0.412000	
Language-COT	0.000519	0.472000	
Shot-Multi Turn	0.000438	0.512000	
Multi Turn-presence_penalty	0.000400	0.512000	
Shot-presence_penalty	0.000375	0.496000	
Question Format-top_p	0.000275	0.604000	
COT-Multi Turn	0.000245	0.612000	
Question Format-max_tokens	0.000226	0.648000	
Question Format-Multi Turn	0.000166	0.708000	
Multi Turn-temperature	0.000135	0.724000	
max_tokens-presence_penalty	0.000128	0.690000	
top_p-presence_penalty	0.000102	0.752000	
Question Paraphrase-Shot	0.000101	0.746000	
COT-top_p	0.000090	0.758000	
Language-presence_penalty	0.000084	0.792000	
Language-max_tokens	0.000073	0.762000	
temperature-presence_penalty	0.000058	0.786000	
Language-Multi Turn	0.000058	0.806000	
Question Paraphrase-top_p	0.000058 0.000049	0.814000 0.816000	
Language-temperature	0.000049	0.818000	
Shot-top_p	0.000037	0.864000	
Question Paraphrase	0.000033	0.860000	
Shot-temperature	0.000033	0.850000	
Language-top_p	0.000029	0.850000	
Question Paraphrase-max_tokens	0.000023	0.802000	
Question Paraphrase-presence_penalty Question Paraphrase-temperature	0.000015	0.904000	
Multi Turn-top_p	0.000015	0.924000	
Shot-max_tokens	0.000013	0.898000	
COT-presence_penalty	0.000013	0.918000	
Question Format-presence_penalty	0.000015	0.946000	
temperature-max_tokens	0.000003	0.970000	
presence_penalty	0.000002	0.960000	
max_tokens-top_p	0.000002	0.970000	
Question Paraphrase-COT	0.000002	0.978000	
Language-Question Paraphrase	0.000000	0.994000	
	0.00000	5.77.000	

## (g) Complete ANOVA results on Mistral Large

# (h) Complete ANOVA results on Mistral Medium

Factor	Effect Size $(\eta^2)$	p-value
Question Format	0.406873	0.000000
Question Format-COT	0.268919	0.000000
COT	0.075852	0.000000
COT-max_tokens	0.026017	0.000000
max_tokens	0.025372	0.000000
Question Format-Multi Turn	0.004796	0.024000
Multi Turn	0.003609	0.058000
Question Format-Shot	0.003338	0.068000
COT-Multi Turn	0.002708	0.100000
Question Format-max_tokens	0.001379	0.230000
Shot	0.001023	0.320000
Language-Question Format	0.001004	0.280000
Multi Turn-max_tokens	0.000881	0.310000
Shot-Multi Turn	0.000830	0.368000
Language-Multi Turn	0.000766	0.372000
Language-COT	0.000765	0.356000
Question Format-Question Paraphrase	0.000644	0.41400
Multi Turn-presence_penalty	0.000456	0.458000
Question Paraphrase	0.000364	0.55000
Question Paraphrase-presence_penalty	0.000353	0.48600
max_tokens-presence_penalty	0.000333	0.62000
	0.000272	0.63800
Shot-top_p Multi Turn-top_p	0.000244	0.57000
	0.000244	0.61800
Language-Question Paraphrase		0.60400
Question Format-presence_penalty	0.000227 0.000227	0.62800
COT-temperature		
Question Paraphrase-Shot	0.000193	0.65800
temperature	0.000185	0.64400
Shot-presence_penalty	0.000178	0.65000
Question Format-temperature	0.000178	0.67200
Question Paraphrase-top_p	0.000163	0.68600
Shot-max_tokens	0.000163	0.68200
top_p	0.000142	0.69000
Question Paraphrase-max_tokens	0.000135	0.68800
Question Paraphrase-Multi Turn	0.000128	0.77200
max_tokens-top_p	0.000110	0.71200
Language	0.000109	0.73600
COT-presence_penalty	0.000087	0.78000
Shot-temperature	0.000087	0.81000
Question Paraphrase-COT	0.000053	0.81800
top_p-presence_penalty	0.000049	0.82400
presence_penalty	0.000038	0.83600
Language-temperature	0.000035	0.84000
Language-top_p	0.000028	0.86400
temperature-max_tokens	0.000017	0.88400
temperature-presence_penalty	0.000015	0.88400
Multi Turn-temperature	0.000007	0.92800
temperature-top_p	0.000006	0.92800
Question Paraphrase-temperature	0.000005	0.95000
COT-top_p	0.000003	0.95200
Language-Shot	0.000003	0.94400
Question Format-top_p	0.000002	0.95200
Language-max_tokens	0.000002	0.95000
Language-presence_penalty	0.000002	0.96000
Shot-COT	0.000002	0.998000

(ii) Complete ANOVA lesur	its on mistrar r	vicuiuiii
Factor	Effect Size $(\eta^2)$	p-value
Question Format	0.430500	0.000000
Question Format-COT	0.274248	0.000000
COT	0.105919	0.000000
COT-max_tokens	0.013038	0.002000
max_tokens	0.012064	0.000000
Question Format-Multi Turn	0.007865	0.004000
Multi Turn	0.005334	0.026000
Shot	0.003097	0.064000
COT-Multi Turn	0.003001	0.090000
Question Format-Shot	0.002782	0.086000
Multi Turn-max_tokens	0.001354	0.252000
Shot-Multi Turn	0.000908	0.336000
Language-Question Paraphrase	0.000891	0.378000
Shot-COT	0.000636	0.434000
Question Paraphrase-top_p	0.000366	0.476000
Question Paraphrase-presence_penalty	0.000312	0.534000
Question Paraphrase-Multi Turn	0.000282	0.578000
Language	0.000254	0.626000
max_tokens-presence_penalty	0.000227	0.650000
Language-Multi Turn	0.000219	0.656000
Multi Turn-presence_penalty	0.000218	0.642000
Question Format-temperature	0.000178	0.682000
Multi Turn-top_p	0.000176	0.672000
Language-COT	0.000163	0.690000
temperature-top_p	0.000148	0.698000
Question Format-max_tokens	0.000118	0.672000
Question Format-Question Paraphrase	0.000120	0.746000
Question Paraphrase-max_tokens	0.000103	0.754000
COT-top_p	0.000091	0.750000
COT-temperature	0.000081	0.764000
Shot-temperature	0.000076	0.806000
top_p	0.000076	0.784000
Question Paraphrase-temperature	0.000076	0.758000
Language-temperature	0.000061	0.798000
Question Format-presence_penalty	0.000053	0.806000
Shot-top_p	0.000048	0.808000
Language-Shot	0.000044	0.832000
Shot-max_tokens	0.000044	0.814000
Question Paraphrase-COT	0.000037	0.854000
presence_penalty	0.000037	0.844000
Question Paraphrase	0.000034	0.850000
Language-top_p	0.000034	0.870000
Language-presence_penalty	0.000024	0.874000
Shot-presence_penalty	0.000022	0.888000
temperature-presence_penalty	0.000019	0.910000
Language-max_tokens	0.000014	0.920000
Language-Question Format	0.000007	0.938000
	0.000007	0.942000
Question Format-top_p	0.000005	0.942000
temperature	0.000003	0.948000
max_tokens-top_p	0.000004	0.956000
Question Paraphrase-Shot	0.000003	0.936000
Multi Turn-temperature		
temperature-max_tokens	0.000001	0.968000
COT-presence_penalty	0.000000	0.986000
top_p-presence_penalty	0.000000	0.994000

#### (i) Complete ANOVA results on Moonshot-v1

Question Format         0.297180         0.00000           Question Format-COT         0.147641         0.00000           Shot-COT         0.130758         0.00000           Shot-COT         0.077565         0.00000           Shot         0.038198         0.00000           COT-Multi Turn         0.018078         0.00000           Language-Shot         0.005572         0.016000           COT-max.tokens         0.002887         0.09000           Language-Multi Turn         0.002887         0.09000           Language-Multi Turn         0.002702         0.110000           Question Format-Multi Turn         0.002700         0.106000           max.tokens         0.002409         0.118000           Multi Turn-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.00066         0.30000           Language-max.tokens         0.000930         0.356000           Shot-max.tokens         0.000930         0.356000           Shot-max.tokens         0.0000930         0.35000           COT-presence-penalty	Factor	Effect Size $(\eta^2)$	p-value
Öuestion Format-COT         0.147641         0.000000           COT         0.130758         0.000000           Shot-COT         0.077565         0.000000           Shot         0.042631         0.000000           Shot         0.038198         0.000000           COT-Multi Turn         0.018078         0.000000           COT-max.tokens         0.002887         0.090000           Language-Multi Turn         0.002824         0.884000           Language-Multi Turn         0.002702         0.110000           Question Format-Multi Turn         0.002709         0.118000           max.tokens         0.002409         0.118000           Language-Question Format         0.002297         0.108000           Multi Turn-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.001677         0.166000           Shot-max.tokens         0.001677         0.166000           Language-max.tokens         0.000966         0.30000           Language-max.tokens         0.000966         0.30000           Language-duestion Paraphrase         0.000961         0.38000           Language-top-pality         0.00042         0.494000           Vopp-resence-penalty <td>Ouestion Format</td> <td>0.297180</td> <td>0.000000</td>	Ouestion Format	0.297180	0.000000
COT         0.130758         0.000000           Shot-COT         0.077565         0.000000           Question Format-Shot         0.042631         0.000000           Shot         0.038198         0.000000           COT-Multi Turn         0.018078         0.000000           Language-Shot         0.005572         0.016000           COT-max.tokens         0.002887         0.090000           Language-Multi Turn         0.002824         0.084000           Question Format-Multi Turn         0.002700         0.106000           max.tokens         0.002409         0.118000           Language-Question Format         0.002297         0.108000           Multi Turn         0.002297         0.108000           Multi Turn-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.00966         0.30000           Language-max.tokens         0.000966         0.300000           Language-max.tokens         0.000930         0.356000           Shot-max.tokens         0.0000930         0.356000           Language-question Paraphrase         0.000793         0.360000           COT-presence-p		0.147641	0.000000
Question Format-Shot         0.042631         0.000000           Shot         0.038198         0.000000           COT-Multi Turn         0.018078         0.000000           Language-Shot         0.005572         0.016000           COT-max.tokens         0.002887         0.090000           Language-Multi Turn         0.002824         0.084000           Language-Multi Turn         0.002700         0.116000           max.tokens         0.002409         0.118000           Language-Question Format         0.002297         0.108000           Multi Turn         0.002185         0.12400           Multi Turn-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.000966         0.300000           Language-max.tokens         0.000930         0.356000           Shot-max.tokens         0.000930         0.356000           Shot-max.tokens         0.0000930         0.356000           Shot-max.tokens         0.0000930         0.36000           COT-presence-penalty         0.000402         0.494000           top-p-presence-penalty         0.00014         0.54900           Multi Turn-temperature         0.000214         0.578000           Question For		0.130758	0.000000
Shot COT-Multi Turn         0.038198         0.000000           Language-Shot         0.005572         0.016000           COT-max_tokens         0.002873         0.09000           Language-Multi Turn         0.002824         0.08400           Language-Wulti Turn         0.002702         0.110000           Question Format-Multi Turn         0.002700         0.106000           max_tokens         0.002409         0.118000           Language-Question Format         0.002297         0.108000           Multi Turn         0.002185         0.124000           Multi Turn-max_tokens         0.001677         0.166000           Question Format-max_tokens         0.000966         0.300000           Language-max_tokens         0.000930         0.356000           Shot-max_tokens         0.000861         0.38000           Shot-max_tokens         0.000793         0.36000           COT-presence_penalty         0.000402         0.494000           toppresence_penalty         0.000414         0.554000           Multi Turn-temperature         0.000274         0.57800           max_tokens-top_p         0.000219         0.662000           guestion Paraphrase-copalty         0.000128         0.724000	Shot-COT	0.077565	0.000000
COT-Multi Turn	Question Format-Shot	0.042631	0.000000
Language-Shot         0,005572         0,016000           COTF-max_tokens         0,002837         0,990000           Language-Multi Turn         0,002824         0,084000           Question Format-Multi Turn         0,002702         0,110000           Multi Turn         0,002709         0,118000           Language-Question Format         0,002297         0,108000           Multi Turn         0,002185         0,124000           Multi Turn-max_tokens         0,001677         0,166000           Question Format-max_tokens         0,000966         0,300000           Language-max_tokens         0,000966         0,300000           Shot-max_tokens         0,000861         0,380000           Shot-max_tokens         0,000861         0,380000           COT-p-resence_penalty         0,000402         0,494000           top-p-presence_penalty         0,00014         0,549000           max_tokens-top_p         0,000219         0,662000           max_tokens-top_p         0,000219         0,662000           guestion Format-top_p         0,000128         0,724000           Question Paraphrase-COT         0,000128         0,724000           Question Paraphrase-COT         0,000128         0,724000 <td>Shot</td> <td>0.038198</td> <td>0.000000</td>	Shot	0.038198	0.000000
COT-max_tokens			
Language-Multi Turn			
Language 0.002702 0.110000 max tokens 0.002409 0.166000 max tokens 0.002409 0.118000 Language-Question Format 0.002297 0.108000 Multi Turn 0.002185 0.124000 Multi Turn-max.tokens 0.001677 0.166000 Question Format-max.tokens 0.001677 0.166000 Language-max.tokens 0.00966 0.300000 Language-max.tokens 0.00966 0.300000 Language-max.tokens 0.009930 0.356000 Language-max.tokens 0.000930 0.356000 Language-max.tokens 0.000931 0.356000 COT-presence.penalty 0.000402 0.494000 topp-presence.penalty 0.00041 0.578000 Multi Turn-temperature 0.000274 0.578000 Multi Turn-temperature 0.000274 0.578000 Multi Turn-temperature 0.000274 0.578000 Shot-presence.penalty 0.000128 0.620000 Shot-presence.penalty 0.000128 0.742000 Question Paraphrase-max.tokens 0.000128 0.742000 Question Paraphrase-max.tokens 0.000128 0.7742000 Question Paraphrase-max.tokens 0.000128 0.756000 Shot-temperature 0.000092 0.786000 Max.tokens-presence.penalty 0.000081 0.776000 Shot-temperature 0.000081 0.776000 Shot-top-p 0.000081 0.776000 Question Paraphrase-Shot 0.000081 0.792000 Question Paraphrase-Shot 0.000081 0.792000 Question Paraphrase-Shot 0.000081 0.782000 Language-temperature 0.000037 0.866000 Question Paraphrase-Multi Turn 0.000033 0.860000 Question Paraphrase-Multi Turn 0.000031 0.782000 Emperature-presence.penalty 0.000031 0.782000 Question Paraphrase-Presence.penalty 0.000031 0.782000 Question Paraphrase-Presence.penalty 0.000031 0.782000 Question Paraphrase-Presence.penalty 0.000031 0.782000 Question Paraphrase-Presence.penalty 0.000031 0.878000 Question Paraphrase-Presence.penalty 0.000031 0.878000 Question Paraphrase 0.000001 0.986000 Question Paraphrase 0.000001 0.986000 Question Paraphrase-temperature 0.00001 0.986000 Question Paraphrase-temperature 0.00001 0.986000 Question Paraphrase-temperature 0.00001 0.996000 Question Paraphrase-temperature 0.00001 0.996000 Question Paraphrase-temperature 0.000001 0.996000 Question Paraphrase-temperature 0.000001 0.996000 Question Paraphrase-temperature 0.000001 0.996000			
Question Format-Multi Turn         0.002700         0.106000           max Jokens         0.002409         0.118000           Language-Question Format         0.002297         0.108000           Multi Turn         0.002185         0.124000           Multi Turn-max Jokens         0.001677         0.166000           Question Format-max Jokens         0.000966         0.300000           Language-Puestion Paraphrase         0.000861         0.380000           Shot-max Jokens         0.000861         0.380000           COT-presence-penalty         0.000402         0.494000           top-p-presence-penalty         0.000414         0.554000           max Jokens-top.p         0.000219         0.562000           max Jokens-top.p         0.000219         0.562000           max Jokens-top.p         0.000128         0.724000           Question Format-top.p         0.000128         0.724000           Question Paraphrase-max.tokens         0.000128         0.724000           Question Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.000021         0.756000           Shot-temperature-top.p         0.000081         0.770000           Shot-top.p         0.000081         0.770			
max.tokens         0.002409         0.118000           Language-Question Format         0.002297         0.108000           Multi Turn         0.002185         0.124000           Multi Turn-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.000966         0.300000           Language-max.tokens         0.000930         0.356000           Shot-max.tokens         0.000793         0.360000           Language-Question Paraphrase         0.000793         0.360000           COTF-presence.penalty         0.000402         0.494000           top.p-presence.penalty         0.000214         0.554000           Multi Turn-temperature         0.000274         0.578000           max.tokens-top.p         0.000219         0.662000           Shot-presence.penalty         0.000203         0.620000           Jop.p         0.000128         0.742000           Question Format-top.p         0.000128         0.724000           Question Paraphrase-max.tokens         0.000128         0.724000           Question Paraphrase-max.tokens         0.000128         0.724000           Question Paraphrase-ce-penalty         0.00015         0.756000           max.tokens-presence-penalty         0.0000			
Language-Question Format         0.002297         0.108000           Multi Turn         0.002185         0.124000           Multi Turn         0.002185         0.124000           Multi Turn-max.tokens         0.001677         0.166000           Question Format-max.tokens         0.000966         0.300000           Language-max.tokens         0.000861         0.380000           Shot-max.tokens         0.000793         0.360000           COT-presence_penalty         0.000402         0.494000           top-p-presence_penalty         0.000214         0.554000           Multi Turn-temperature         0.000274         0.578000           max.tokens-top-p         0.000219         0.662000           Shot-presence-penalty         0.000219         0.662000           Shot-presence-penalty         0.000128         0.724000           Question Format-top-p         0.000128         0.724000           Question Paraphrase-eCT         0.000128         0.724000           Question Paraphrase-eCT         0.000128         0.724000           Shot-temperature         0.000081         0.776000           Shot-temperature-top-p         0.000081         0.776000           Shot-top-p         0.000081         0.776000			
Mulfi Turn         0.002185         0.124000           Mulfi Turn-max Jokens         0.001677         0.166000           Question Format-max Jokens         0.000966         0.300000           Language-max Jokens         0.000930         0.356000           Shot-max Jokens         0.000861         0.380000           Language-Question Paraphrase         0.000793         0.360000           COT-presence-penalty         0.000412         0.494000           top-p-presence-penalty         0.000214         0.554000           max Jokens-top.p         0.000219         0.662000           shot-presence-penalty         0.000219         0.662000           max Jokens-top.p         0.000128         0.738000           shot-presence-penalty         0.000128         0.742000           Question Format-top.p         0.000128         0.742000           Question Paraphrase-COT         0.00115         0.756000           Shot-temperature         0.00092         0.786000           max Jokens-presence-penalty         0.00092         0.786000           max Jokens-presence-penalty         0.00081         0.770000           shot-top.p         0.000081         0.776000           shot-top.p         0.000081         0.776000 <td></td> <td></td> <td></td>			
Multi Turn-max_tokens         0.001677         0.166000           Question Format-max_tokens         0.000930         0.356000           Language-max_tokens         0.000930         0.356000           Shot-max_tokens         0.000861         0.380000           COF-presence_penalty         0.000402         0.494000           COF-presence_penalty         0.00041         0.554000           Multi Turn-temperature         0.000274         0.578000           max_tokens-top_p         0.000219         0.662000           Shot-presence_penalty         0.000219         0.662000           Shot-presence_penalty         0.000213         0.520000           top_p         0.000219         0.662000           Shot-presence_penalty         0.000128         0.724000           Question Format-top_p         0.000128         0.724000           Question Paraphrase-eCT         0.000128         0.724000           Question Paraphrase-epenalty         0.000128         0.724000           Shot-temperature         0.000081         0.776000           Shot-top_p         0.000081         0.776000           Shot-top_p         0.000081         0.776000           Question Paraphrase-Multi Turn         0.000081         0.782000			
Question Format-max_tokens         0.000966         0.300000           Language-max_tokens         0.000861         0.386000           Shot-max_tokens         0.000861         0.380000           COT-presence_penalty         0.000402         0.494000           top_p-presence_penalty         0.00014         0.554000           Multi Turn-temperature         0.000274         0.578800           max_tokens-top_p         0.000219         0.662000           Shot-presence_penalty         0.000128         0.20000           top_p         0.000142         0.694000           Question Format-top_p         0.000128         0.724000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-COT         0.000128         0.724000           Question Paraphrase-CoT         0.00015         0.756000           Shot-temperature         0.000081         0.776000           Shot-temperature-top_p         0.00081         0.776000           Max_tokens-presence-penalty         0.00081         0.776000           Shot-temperature-top_p         0.00081         0.776000           Language-temperature         0.00081         0.782000           temperature-presence-penalty         0.000081 <td></td> <td></td> <td></td>			
Language-max_tokens         0.000930         0.356000           Shot-max_tokens         0.000861         0.380000           Language-Question Paraphrase         0.000793         0.360000           COT-presence_penalty         0.000402         0.494000           top-p-presence_penalty         0.000314         0.578000           Multi Turn-temperature         0.000274         0.578000           max_tokens-top-p         0.000219         0.662000           Shot-presence_penalty         0.000203         0.620000           top-p         0.000124         0.694000           Question Format-top-p         0.000128         0.742000           Question Format-top-p         0.000128         0.724000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.000081         0.776000           Shot-temperature-op-p         0.000081         0.779000           Shot-top-p         0.000081         0.776000           Question Paraphrase-Shot         0.00081         0.782000           Language-temperature         0.000081         0.782000           temperature-presence-penalty         0.000033			
Shot-max_tokens			
Language-Question Paraphrase         0.000793         0.360000           COT-presence-penalty         0.000402         0.494000           top-p-presence-penalty         0.000314         0.554000           Multi Turn-temperature         0.000274         0.578000           max_tokens-top-p         0.000219         0.662000           Shot-presence-penalty         0.000203         0.620000           top-p         0.000128         0.742000           Question Format-top-p         0.000128         0.742000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-cOT         0.000115         0.756000           Max_tokens-presence-penalty         0.000081         0.779000           Shot-temperature         0.000081         0.779000           Shot-top-p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.782000           Language-temperature         0.000081         0.782000           Language-temperature         0.000081         0.782000           Lemperature-presence-penalty         0.000053         0.858000           Shot-Multi Turn         0.000			
COT-presence_penalty			
top_p-presence_penalty         0.000314         0.554000           Mulii Turn-temperature         0.000274         0.578000           max_tokens-top_p         0.000219         0.662000           Shot-presence_penalty         0.000219         0.662000           Shot-presence_penalty         0.000128         0.742000           Question Format-top_p         0.000128         0.724000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-ency         0.000115         0.756000           Question Paraphrase-COT         0.00015         0.756000           max_tokens-presence_penalty         0.00081         0.770000           Shot-temperature         0.000081         0.770000           Shot-top_p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.782000           Language-temperature         0.000081         0.782000           temperature-presence_penalty         0.000053         0.858000           Shot-Multi Turn         0.00033         0.858000           Shot-Multi Turn         0.00037         0.856000           Question Format-temperature         0			
Multi Turn-temperature         0.000274         0.578000           max_tokens-top_p         0.000219         0.662000           Shot-presence_penalty         0.000203         0.620000           top_p         0.000142         0.694000           Question Format-top_p         0.000128         0.724000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.000092         0.786000           max_tokens-presence_penalty         0.000081         0.770000           Shot-top_p         0.000081         0.779000           Shot-top_p         0.000081         0.776000           Question Paraphrase-Shot         0.00081         0.776000           Question Paraphrase-Shot         0.00081         0.776000           Language-temperature         0.000081         0.782000           Question Paraphrase-Multi Turn         0.000033         0.86000           Question Paraphrase-Multi Turn         0.000033         0.856000           Question Format-temperature         0.000037         0.856000           Question Paraphrase-presence-penalty         0.000037         0.872000           Question Format-presence-penalty			
max_tokens-top_p         0.000219         0.662000           Shot-presence_penalty         0.00023         0.52000           top_p         0.000142         0.694000           Question Format-top_p         0.000128         0.742000           Question Paraphrase-emax_tokens         0.000128         0.724000           Question Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.000092         0.786000           max_tokens-presence_penalty         0.000081         0.770000           Shot-top_p         0.000081         0.770000           Question Paraphrase-Shot         0.00081         0.776000           Question Paraphrase-Shot         0.000081         0.782000           Language-temperature         0.000081         0.782000           temperature-presence_penalty         0.000053         0.858000           Question Paraphrase-Multi Turn         0.000033         0.858000           Shot-Multi Turn         0.000037         0.856000           COT-top_p         0.000037         0.856000           COT-top_p         0.000037         0.856000           COT-top_p         0.000037         0.872000           Question Paraphrase-presence_penalty         0.000037         0.8720			
Shot-presence_penalty			
top_p*         0.000142         0.694000           Question Format-top_p         0.000128         0.742000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.00092         0.786000           max_tokens-presence-penalty         0.000081         0.770000           Shot-top_p         0.000081         0.7792000           temperature-top_p         0.000081         0.804000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence-penalty         0.000053         0.86000           Question Paraphrase-Multi Turn         0.000033         0.858000           Question Format-temperature         0.000037         0.856000           COT-top_p         0.000037         0.866000           Question Paraphrase-presence-penalty         0.000037         0.878000           Question Paraphrase         0.000037         0.878000           Question Paraphrase         0.000037         0.878000           Question Paraphrase         0.000037         0.878000           Question Paraphrase         0.00			
Question Format-top_p         0.000128         0.742000           Question Paraphrase-max_tokens         0.000128         0.724000           Question Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.000092         0.786000           max_tokens-presence-penalty         0.000081         0.770000           Shot-top_p         0.000081         0.776000           Shot-top_p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence_penalty         0.000053         0.860000           Question Paraphrase-Multi Turn         0.000033         0.858000           Question Format-temperature         0.000037         0.856000           Question Paraphrase-presence_penalty         0.000037         0.856000           Question Paraphrase-presence_penalty         0.000037         0.872000           Question Format-presence_penalty         0.000037         0.872000           Question Format-presence_penalty         0.000030         0.878000           Question Paraphrase         0.000024         0.864000           Language-presence_penalty         0.000019         0.884000 </td <td></td> <td></td> <td></td>			
Question Paraphrase-max.tokens         0.000128         0.724000           Question Paraphrase-COT         0.00015         0.756000           Shot-temperature         0.000092         0.786000           Shot-temperature         0.000081         0.770000           Shot-top_p         0.000081         0.792000           temperature-top_p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence-penalty         0.000053         0.860000           Question Paraphrase-Multi Turn         0.000053         0.858000           Shot-Multi Turn         0.00037         0.856000           Question Format-temperature         0.000037         0.856000           Question Paraphrase-presence-penalty         0.00037         0.856000           Question Paraphrase-presence-penalty         0.000037         0.872000           Question Pormat-presence-penalty         0.00037         0.878000           Question Pormat-presence-penalty         0.00037         0.878000           Question Paraphrase         0.000024         0.864000           Language-presence-penalty         0.000019         0.984000			
Öuestion Paraphrase-COT         0.000115         0.756000           Shot-temperature         0.00092         0.786000           max.tokens-presence.penalty         0.000081         0.770000           Shot-top.p         0.000081         0.776000           Generature-top.p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence.penalty         0.000053         0.858000           Question Paraphrase-Multi Turn         0.00053         0.858000           Shot-Multi Turn         0.00037         0.856000           Question Format-temperature         0.00037         0.856000           Question Paraphrase-presence-penalty         0.00037         0.882000           Question Paraphrase-presence-penalty         0.00037         0.872000           Question Format-presence-penalty         0.00037         0.872000           Question Paraphrase         0.000037         0.878000           Question Paraphrase         0.000034         0.878000           Language-presence-penalty         0.000019         0.884000           Language-COT         0.000019         0.884000           Language-top-			
Shot-temperature         0.00092         0.786000           max.tokens-presence_penalty         0.000081         0.770000           Shot-top_p         0.000081         0.772000           Images of the perature top_p         0.000081         0.776000           Question Paraphrase-Shot         0.00081         0.804000           Language-temperature         0.000081         0.804000           temperature-presence_penalty         0.000053         0.860000           Question Paraphrase-Multi Turn         0.000053         0.858000           Shot-Multi Turn         0.000037         0.856000           Question Format-temperature         0.000037         0.856000           Question Paraphrase-presence-penalty         0.000037         0.872000           Question Paraphrase         0.00037         0.872000           Question Pormat-presence-penalty         0.00037         0.872000           Question Permat-presence-penalty         0.000030         0.878000           Question Paraphrase         0.000024         0.864000           Language-presence-penalty         0.000019         0.884000           Language-presence-penalty         0.000019         0.884000           Umbit Turn-top-p         0.00001         0.958000           <			
max.tokens-presence_penalty         0.000081         0.770000           Shot-top_p         0.000081         0.792000           temperature-top_p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence_penalty         0.000053         0.858000           Question Paraphrase-Multi Turn         0.000045         0.812000           Question Format-temperature         0.000037         0.856000           COT-top_p         0.00037         0.856000           Question Paraphrase-presence_penalty         0.00037         0.872000           Question Format-presence_penalty         0.00037         0.872000           Question Format-presence_penalty         0.00033         0.878000           Question Format-presence_penalty         0.000034         0.864000           Language-presence_penalty         0.000019         0.884000           Language-COT         0.000019         0.884000           Multi Turn-top-p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Fo			
Shot-top-p         0.000081         0.792000           temperature-top.p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence-penalty         0.000083         0.860000           Question Paraphrase-Multi Turn         0.000033         0.858000           Shot-Multi Turn         0.000037         0.856000           Question Format-temperature         0.000037         0.856000           OU-stion Paraphrase-presence-penalty         0.000037         0.872000           Question Paraphrase-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000037         0.872000           Question Paraphrase         0.000024         0.864000           Language-presence-penalty         0.000019         0.884000           Language-presence-penalty         0.000019         0.884000           Multi Turn-top-p         0.000019         0.884000           Multi Turn-top-p         0.00001         0.958000           Question Format-Question Paraphrase         0.000002         0.958000           Question Format-Question Paraphrase         0.000001         0.970000 <td></td> <td></td> <td></td>			
temperature-top_p         0.000081         0.776000           Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           Question Paraphrase-Multi Turn         0.000053         0.858000           Question Paraphrase-Multi Turn         0.000045         0.812000           Question Format-temperature         0.000037         0.856000           COT-top_p         0.000037         0.866000           Question Paraphrase-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000033         0.878000           Question Format-presence-penalty         0.000024         0.864000           Language-Presence-penalty         0.000019         0.904000           Language-COT         0.000019         0.884000           Multi Turn-top_p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000001         0.966000 </td <td></td> <td></td> <td></td>			
Question Paraphrase-Shot         0.000081         0.804000           Language-temperature         0.000081         0.782000           temperature-presence-penalty         0.000053         0.860000           Question Paraphrase-Multi Turn         0.000033         0.858000           Shot-Multi Turn         0.000037         0.856000           Question Format-temperature         0.000037         0.856000           COT-top-p         0.000037         0.882000           Question Paraphrase-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000030         0.878000           Question Paraphrase         0.000024         0.864000           Language-presence-penalty         0.000019         0.884000           Language-presence-penalty         0.000019         0.884000           Language-presence-penalty         0.000019         0.884000           Language-presence-penalty         0.000019         0.884000           Language-presence-penalty         0.000019         0.958000           Usestion Paraphrase         0.000002         0.958000           Question Format-question Paraphrase         0.000001         0.966000           Question Paraphrase-temperature         0.000001         0.970000			
Language-temperature         0.000081         0.782000           temperature-presence_penalty         0.00053         0.86000           Question Paraphrase-Multi Turn         0.000045         0.812000           Shot-Multi Turn         0.000045         0.812000           Question Format-temperature         0.00037         0.856000           COT-top-p         0.000037         0.886000           Question Paraphrase-presence_penalty         0.00037         0.872000           Multi Turn-presence_penalty         0.00037         0.872000           Question Format-presence_penalty         0.000030         0.878000           Question Format-presence_penalty         0.000024         0.864000           Language-presence_penalty         0.000019         0.94000           Language-COT         0.000019         0.884000           Multi Turn-top-p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000001         0.966000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.976000           temp			
temperature-presence_penalty 0.000053 0.850000 Ouestion Paraphrase-Multi Turn 0.000053 0.858000 Shot-Multi Turn 0.000053 0.858000 Shot-Multi Turn 0.000045 0.812000 Question Format-temperature 0.000037 0.856000 COT-top_p properature 0.000037 0.866000 Question Paraphrase-presence_penalty 0.000037 0.872000 Question Paraphrase 0.000037 0.872000 Question Paraphrase 0.000030 0.878000 Question Paraphrase 0.000030 0.878000 Question Paraphrase 0.000030 0.878000 Question Paraphrase 0.000019 0.904000 Language-presence_penalty 0.000019 0.904000 Uanguage-COT 0.000019 0.884000 0.00019 0.904000 Uanguage-COT 0.000019 0.904000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.9050000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.9050000 0.905000 0.905000 0.905000 0.905000 0.905000 0.905000 0.90500			
Question Paraphrase-Multi Turn         0.000053         0.858000           Shot-Multi Turn         0.000045         0.812000           Question Format-temperature         0.000037         0.856000           COT-top-p         0.00037         0.866000           Question Paraphrase-presence-penalty         0.000037         0.872000           Multi Turn-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000030         0.8788000           Question Paraphrase         0.000024         0.864000           Language-presence-penalty         0.000019         0.94000           Language-presence-penalty         0.000019         0.884000           Multi Turn-top-p         0.000014         0.922000           temperature-max.tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000001         0.966000           Question Paraphrase-temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence-penalty         0.000001         0.968000			
Shot-Multi Turn         0.000045         0.812000           Question Format-temperature         0.000037         0.856000           COT-top-D         0.000037         0.866000           Question Paraphrase-presence-penalty         0.000037         0.872000           Multi Turn-presence-penalty         0.000037         0.872000           Question Format-presence-penalty         0.000030         0.878000           Question Paraphrase         0.00024         0.864000           Language-presence-penalty         0.000019         0.904000           Language-COT         0.000019         0.884000           Multi Turn-top-p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence-penalty         0.000001         0.992000			
Question Format-temperature         0.000037         0.856000           COT-top-p         0.00037         0.866000           Question Paraphrase-presence penalty         0.000037         0.872000           Multi Turn-presence penalty         0.000037         0.872000           Question Format-presence-penalty         0.000030         0.878800           Question Paraphrase         0.000024         0.864000           Language-presence-penalty         0.000019         0.94000           Language-COT         0.000019         0.884000           Multi Turn-top-p         0.00001         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.954000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence-penalty         0.000000         0.992000			
ČOT-top.p         0.000037         0.866000           Question Paraphrase-presence_penalty         0.000037         0.882000           Multi Turn-presence_penalty         0.000037         0.872000           Question Format-presence_penalty         0.000030         0.878000           Question Paraphrase         0.00024         0.864000           Language-presence_penalty         0.000019         0.94000           Language-COT         0.00019         0.884000           Multi Turn-top.p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top.p         0.000001         0.968000           presence_penalty         0.000001         0.998000			
Question Paraphrase-presence penalty         0.000037         0.882000           Multi Turn-presence_penalty         0.000030         0.8782000           Question Format-presence_penalty         0.000024         0.864000           Language-presence_penalty         0.000019         0.904000           Language-COT         0.000019         0.884000           Multi Turn-top-p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COTI-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence_penalty         0.000001         0.992000			
Multi Turn-presence_penalty         0.000037         0.872000           Question Format-presence_penalty         0.000030         0.878000           Question Paraphrase         0.000024         0.864000           Language-presence_penalty         0.000019         0.904000           Language-COT         0.000019         0.884000           Multi Turn-top_p         0.000014         0.922000           temperature-max_tokens         0.00002         0.958000           COT-temperature         0.000002         0.954000           Question Format-Question Paraphrase         0.000001         0.966000           question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top_p         0.000001         0.968000           presence_penalty         0.000001         0.992000			
Question Format-presence penalty         0.000030         0.878000           Question Paraphrase         0.000024         0.864000           Language-presence penalty         0.000019         0.904000           Language-COT         0.000019         0.884000           Multi Turn-top-p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence-penalty         0.000001         0.992000			
Question Paraphrase         0.000024         0.864000           Language-presence penalty         0.00019         0.94000           Language-COT         0.000019         0.884000           Multi Turn-top_p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.958000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top_p         0.000001         0.968000           presence_penalty         0.000000         0.992000			
Language-presence penalty         0.000019         0.904000           Language-COT         0.000019         0.884400           Multi Turn-top.p         0.000014         0.922000           temperature-max_tokens         0.000002         0.958000           COT-temperature         0.000002         0.954000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top.p         0.000001         0.968000           presence-penalty         0.000000         0.992000			
Language-COT         0.000019         0.884000           Multi Turn-top-p         0.00014         0.922000           temperature-max tokens         0.000002         0.958000           COT-temperature         0.000002         0.954000           Question Format-Question Paraphrase         0.000001         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence-penalty         0.000000         0.992000			
Multi Turn-top_p         0.000014         0.922000           temperature-max_tokens         0.00002         0.958000           COT-temperature         0.00002         0.954000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top_p         0.000001         0.968000           presence_penalty         0.000001         0.992000			
temperature-max tokens         0.000002         0.958000           COTI-temperature         0.00002         0.954000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top.p         0.000001         0.968000           presence.penalty         0.000000         0.992000			
COT-temperature         0.00002         0.954000           Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top.p         0.000001         0.968000           presence_penalty         0.000000         0.992000			
Question Format-Question Paraphrase         0.000002         0.968000           Question Paraphrase-temperature         0.00001         0.966000           temperature         0.000001         0.970000           Language-top.p         0.000001         0.968000           presence.penalty         0.000000         0.992000			
Question Paraphrase-temperature         0.000001         0.966000           temperature         0.000001         0.970000           Language-top-p         0.000001         0.968000           presence-penalty         0.000000         0.992000			
temperature         0.000001         0.970000           Language-top-p         0.00001         0.96800           presence-penalty         0.00000         0.992000			
Language-top_p         0.000001         0.968000           presence_penalty         0.00000         0.992000			
presence_penalty 0.000000 0.992000			
Question Paraphrase-top_p 0.000000 0.998000	Question Paraphrase-top_p	0.000000	0.998000

# A.5 EXPLORING THE POTENTIAL OF LLMs IN SAFETY-CRITICAL APPLICATIONS

In this section, we analyze the relationships between the latency and accuracy of LLMs, as well as between latency and hardware architectures, based on online evaluation. On one hand, we conduct online testing to assess the accuracy of LLMs under different configuration spaces in terms of the "tail to quality" (Yang et al., 2022) metric. Here, "tail to quality" refers to the ratio of the number of tasks correctly completed within a specified threshold to the total number of tasks. Figure 8a illustrates the performance of various LLMs under the "Tail to Quality" metric, showing how their quality scores evolve across different threshold values. Among the models, deepseek (green curve) consistently demonstrates the highest quality across all thresholds, outperforming the others. Doubao (blue curve) and qwen (gray curve) follow, with doubao approaching deepseek's performance at higher thresholds. Kimi (brown curve) and qwen25 (cyan curve) exhibit relatively lower quality, though qwen25 shows rapid improvement at lower thresholds before plateauing. Overall, the chart highlights deepseek's superior capability in handling tail data, while qwen25's growth in quality becomes limited at higher thresholds.

On the other hand, following a similar approach as for accuracy, we obtain the 95% and 99% confidence intervals for latency, as shown in Figure 8b. It can be seen that, for most models, latency and accuracy on AIME'2024 are positively correlated. Notably, doubao-1.5-pro and qwen2.5 achieve relatively low latency while maintaining high accuracy. In contrast, gpt-4.1 and qwen-plus exhibit the opposite trend: they achieve lower accuracy despite higher latency.

