

BENCHMARKING LARGE LANGUAGE MODEL BENCHMARKS: POPULAR BENCHMARKS VS. HUMAN PERCEPTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Benchmarks play a critical role as a measure of large language model (LLM) capabilities. However, whether LLM performance on benchmarks is similar to their real-world performance, especially human perception of their outputs, remains questionable. This study specifically focuses on whether **LLM performance on benchmarks is similar to human perception**. The study investigates this gap by quantifying the similarity between LLM rankings derived from benchmarks and LLM rankings generated from human votes on the prominent LMArena platform. It systematically compares benchmark rankings against rankings in corresponding task-specific categories in LMArena for over 100 top-tier LLMs. The findings reveal that LLM performance on several popular benchmarks has low similarity with human perception, even though these benchmarks are more recent or challenging. The results highlight limitations in current benchmarking practices and underscore the need for evaluation frameworks that more accurately reflect the human perception and real-world performance of LLMs.

1 INTRODUCTION

Benchmarking has become a cornerstone for evaluating large language models (LLMs). When new LLMs are released, their capabilities are typically assessed and compared using standardized benchmarks, such as (Hendrycks et al., 2021b;a), LiveCodeBench (Jain et al., 2025), and Humanity’s Last Exam (Phan et al., 2025). These benchmarks offer a convenient framework for quantifying LLM performance. In recent months, both open-source models and closed-source models have reported state-of-the-art (SOTA) or near-SOTA results on major benchmarks. Examples of these open-source models include the Qwen 3 series (Yang et al., 2025), DeepSeek V3.1 (DeepSeek, 2025), and Kimi-K2 (Team et al., 2025a). Examples of these closed-source models include GPT-5 (OpenAI, 2025), Gemini 2.5 Pro (Comanici et al., 2025), and Claude Opus 4.1 (Anthropic, 2025).

However, as Sam Altman has observed, users increasingly care *less* about benchmark scores and *more* about who is using an LLM and what value they derive from it (Theo Von, 2025). This suggests that despite the rapid progress reflected in SOTA benchmark results, the gap between human perception of LLM outputs and benchmark performance is widening, and trust in benchmark results is declining (Pandey, 2025). These issues are often attributable to cases where AI systems that appear impressive in benchmark evaluations perform poorly in real-world applications (Petrosino, 2025). A likely reason is that models, researchers, and vendors often optimize for the benchmark rather than the real task (Liubimov, 2025). The above observations lead to our central research question:

How does LLM performance on benchmarks differ from human perception?

To explore this, we searched for a benchmark that is fully grounded in real human perception. We identified *LMArena* (Zheng et al., 2023), an open platform where users interact with and compare leading LLMs. By allowing side-by-side comparisons and collecting votes for the better response, the platform enables the community to help rank over 250 LLMs. Additionally, LMArena assesses LLMs’ strengths and weaknesses in a more granular way by grouping user-submitted prompts into

task-specific categories, such as math prompts and coding prompts. LMArena provides a ranking of LLMs for each category of prompts.

Benchmarks can also be classified according to the skill they measure, like mathematical reasoning or coding ability. Given this categorization, our approach is centered on a key assumption: If a benchmark’s ranking of LLMs differs from the ranking of LLMs on the corresponding category in LMArena, then LLM performance on that benchmark must differ from human perception. Guided by this assumption, we compare each benchmark’s ranking against rankings in corresponding task-specific categories in LMArena.

We draw the following conclusions from the above comparison:

- We found that some popular benchmarks have rankings not similar to LMArena’s, like *Humanity’s Last Exam*, *FACTS Grounding*, and *IFEval*. LLM performance on these benchmarks may differ from human perception.
- LLM performance on neither more difficult nor easier benchmarks is necessarily more similar to human perception.
- LLM performance on newer benchmarks is not necessarily more similar to human perception, except for benchmarks that evaluate instruction following or the creative writing ability of LLMs. These benchmarks are being improved to better resemble human perception of LLM outputs.
- Compared with other abilities (e.g., coding ability and instruction following ability), the math ability of LLMs can be relatively better evaluated by current benchmarks.

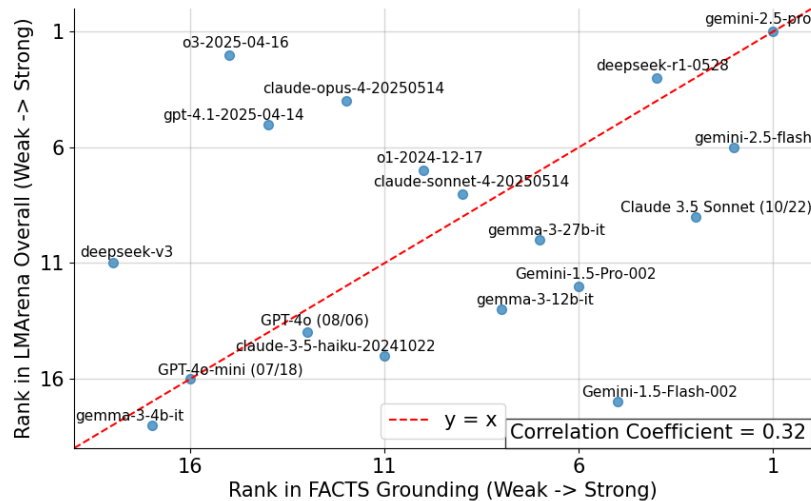


Figure 1: Comparison of large language model (LLM) ranking in FACTS Grounding and the overall ranking in LMArena. Each point (x_i, y_i) represents an LLM ranked x_i -th in A' and y_i -th in B' , where A' and B' are adjusted rankings based only on LLMs that appear in both original rankings A and B . In this case, A and B represent the ranking in FACTS Grounding and the overall ranking in LMArena, respectively. The Spearman rank-order correlation coefficient of these two rankings is merely 0.32. We can conclude that LLM performance on FACTS Grounding may differ from human perception.

These findings expose both the strengths and limitations of popular benchmarks and underscore the urgent need for more representative benchmarks—especially for open-ended, creative, and user-centric tasks.

2 RELATED WORK

Among the benchmarks we investigated, some use LLMs as judges. Most recent work on benchmarking focused on the reliability of LLMs as judges. Zheng et al. (2023) found that powerful LLM judges like GPT-4 can achieve over 80% agreement with human evaluations. However, Wang et al. (2024a) revealed systematic biases when using LLMs as judges, indicating flaws in this evaluation paradigm. Furthermore, Koo et al. (2024) thoroughly investigated the correlation between human preferences and machine preferences, calculating an average Rank-Biased Overlap (RBO) score of just 44%, indicating significant deviation between machine and human preferences. Li et al. (2025a) noted that current LLM evaluators still have certain limitations. However, they also found that when LLMs carefully consider various criteria before giving overall scores, they can achieve higher correlation with human assessments. Nonetheless, none of the aforementioned work has focused on whether benchmark results that do not use LLMs as judges are similar to human perception.

Before LLMs became prevalent, researchers studied the reliability of benchmarks that do not use LLMs as judges, but these benchmarks were not used to evaluate LLMs either. Kiela et al. (2021) pointed out that while contemporary models perform well on benchmarks, they still underperform on simple challenge cases and real-world scenarios, reflecting a disconnect between benchmark performance and practical needs. Dehghani et al. (2021) revealed significant variations in model rankings across different benchmarks. For example, the average Kendall’s rank correlation between subtasks and composite scores in the SuperGLUE benchmark was only 0.648, suggesting that even different tests within the same benchmark may produce inconsistent evaluations, and the correlation between subtasks and average results is imperfect. However, most of the benchmarks focus on performance averaged over many tasks, and the question of how to reliably evaluate and tune models trained for individual tasks in this regime has not been addressed (Shimabucoro et al., 2024). Meanwhile, Liang et al. (2023) studied the correlation between model rankings on accuracy, robustness, and fairness. Yet, they did not explore the relationship between accuracy and human perception of model outputs.

3 METHOD

We quantify the similarity between LLM rankings derived from benchmarks and rankings collected through LMArena by calculating the correlation between them. This section addresses the following questions:

1. Which LLMs do we investigate?
2. Which rankings from benchmarks do we investigate?
3. Which rankings from the LMArena platform do we investigate?
4. How do we calculate the correlation between two rankings?

3.1 INVESTIGATED SET OF LLMs

Since powerful LLMs with larger parameter sizes often receive more attention and are most frequently employed for human-facing services, we focus our investigation specifically on the set of LLMs that have scores not lower than 1300 on the overall LMArena leaderboard. This set of LLMs has a size of 114, and in other words, we investigate the LLMs ranked top 114 on the overall LMArena leaderboard.

3.2 INVESTIGATED LLM RANKINGS FROM BENCHMARKS

We choose 24 public benchmarks on the basis of the wide adoption by top-tier LLMs today, and collect the rankings of the above LLMs from benchmark leaderboards.

We organize our benchmarks into four categories: Question Answering (QA), Mathematics, Code, and Alignment. This classification reflects the core skills emphasized in evaluations for state-of-the-art models such as GPT-5 (OpenAI, 2025), Gemini 2.5 Pro (Comanici et al., 2025), Qwen 3 (Yang et al., 2025), and DeepSeek-R1 (DeepSeek-AI et al., 2025), as well as the organizational structure of widely recognized leaderboards like the Open LLM Leaderboard (Aidar Myrzakhan, 2024). We also draw on insights from a recent comprehensive survey on LLM evaluation (Cao et al., 2025).

Table 1: Investigated benchmarks and their categorization.

Question Answering	Mathematics	Code	Alignment
GPQA	MGSM	HumanEval	IFEval
MMLU-Pro	MATH-500	LiveCodeBench	Arena-Hard
SimpleQA	FrontierMath	SWE-Bench Verified	WritingBench
FACTS Grounding	AIME	Aider Polyglot	Creative Writing v3
Humanity’s Last Exam	HMMT February 2025	Terminal-Bench	IFBench
SuperGPQA		SciCode	
ARC-AGI-2		IOI	

- Question Answering:** Evaluates question answering ability, including factual knowledge recall and logical reasoning. This category includes 7 benchmarks: *GPQA* (Rein et al., 2024), *MMLU-Pro* (Wang et al., 2024b), *SimpleQA* (Wei et al., 2024), *FACTS Grounding* (Jacovi et al., 2025), *Humanity’s Last Exam* (Phan et al., 2025), *SuperGPQA* (Team et al., 2025b), *ARC-AGI-2* (Chollet et al., 2025).
- Mathematics:** Evaluates mathematical problem-solving ability. This category includes 5 benchmarks: *MGSM* (Shi et al., 2023), *MATH-500* (Lightman et al., 2024), *FrontierMath (Tier 1-3)* (Glazer et al., 2025), *AIME* (AoPS, 2025), *HMMT February 2025* (Harvard–MIT Mathematics Tournament (HMMT), 2025).
- Code:** Evaluates code generation and understanding ability. This category includes 7 benchmarks: *HumanEval* (Chen et al., 2021), *LiveCodeBench* (Jain et al., 2025), *SWE-bench Verified* (Jimenez et al., 2024; Chowdhury et al., 2024), *Aider Polyglot* (Aider, 2024), *Terminal-Bench* (Team, 2025), *SciCode* (Tian et al., 2024), *IOI* (International Olympiad in Informatics, 2025).
- Alignment:** Evaluates instruction following and creative writing ability. This category includes 5 benchmarks: *IFEval* (Zhou et al., 2023), *IFBench* (Pyatkin et al., 2025), *ArenaHard* (Li et al., 2025b), *WritingBench* (Wu et al., 2025), *Creative Writing v3* (Paech, 2025).

Table 1 summarizes the benchmarks within each of our four categories. Descriptions of these benchmarks and the sources of the ranking data on these benchmarks can be found in the appendix, in Table 7. Table 9 shows the benchmarks we investigated but not adopted, and the reasons can also be seen in the appendix.

3.3 INVESTIGATED LLM RANKINGS FROM LMARENA

Correspondingly, besides the overall ranking in LMArena, LMArena groups human prompts into various categories, and for each category, it generates an LLM ranking. Each ranking reflects human perception of how well LLMs respond to prompts within this category. To compare the benchmark rankings more rigorously, in addition to the **overall ranking in LMArena**, we investigate three other categories in LMArena, each with its own ranking: **Coding**, **Math**, and **Instruction Following**. More details can be found in the LMArena Blog (<https://news.lmarena.ai/chatbot-arena-categories>).

We collect the scores and rankings of LLMs in the set of LLMs defined in Section 3.1 on LMArena.

3.4 CORRELATION COMPUTATION

Given that our rankings are derived from benchmark scores with heterogeneous distributions, we compute the Spearman rank-order correlation coefficient (Virtanen et al., 2025) between rankings. The Spearman correlation is a non-parametric statistic that assesses the strength and direction of association between two variables without requiring assumptions of linearity or normality.¹

A ranking is an ordered list of LLMs sorted by score or accuracy. Suppose we have two rankings:

$$A = [a_1, a_2, a_3, \dots, a_n], \quad B = [b_1, b_2, b_3, \dots, b_m],$$

¹We do not use the Pearson correlation coefficient, which is a parametric measure that presumes approximate normality of the variables involved.

where each element represents an LLM.

We then remove all LLMs in A that do not appear in B to obtain A' . Similarly, we remove all LLMs in B that do not appear in A to obtain B' . The resulting lists A' and B' have the same length and contain the same set of LLMs, denoted as $A \cap B$. The number of LLMs in $A \cap B$ is defined as $N \triangleq |A \cap B|$, representing the number of data points involved in computing the correlation coefficient.

In a Cartesian coordinate system (refer to Figure 1), each LLM in $A \cap B$ is plotted as a point (x_k, y_k) , where x_k is the rank of that LLM in A' and y_k is the rank in B' . Using these N pairs $(x_1, y_1), (x_2, y_2), \dots$, we compute the Spearman correlation coefficient ρ as

$$\rho = 1 - \frac{6 \sum_{k=1}^N (x_k - y_k)^2}{N(N^2 - 1)}.$$

In practice, we use the library function `scipy.stats.spearmanr` from the Python library `scipy` to calculate ρ . For significance (p-value), when $N < 30$, we apply a permutation test with `scipy.stats.permutation_test` and `permutation_type='pairings'`, which permutes one ranking under the null hypothesis of independence and recalculates the correlation. When $N \geq 30$, the p-value returned by `scipy.stats.spearmanr` is used as a sufficiently accurate approximation.

In Section 3.3, we collect 4 distinct rankings from LMArena and its categories. Let $\{A_1, A_2, A_3, A_4\}$ denote these rankings, and $\{B_1, B_2, \dots, B_{24}\}$ denote the rankings from the 24 benchmarks. For each pair (A_i, B_j) in the Cartesian product of the above two sets, we compute $\rho(A_i, B_j)$, as shown in the schematic table (Table 2): rows represent B_1 – B_{24} , columns represent A_1 – A_4 , and we compute $\rho(A_i, B_j)$ in each cell. All correlation coefficients are reported along with their p-values.

Table 2: A schematic table showing Spearman correlation coefficients $\rho(A_i, B_j)$ calculated between different task-specific rankings in LMArena $\{A_1, A_2, A_3, A_4\}$ and twenty-four benchmark rankings $\{B_1, B_2, \dots, B_{24}\}$. For the computed value of each $\rho(A_i, B_j)$, please refer to Table 5 in the appendix.

	A_1	A_2	A_3	A_4
B_1	$\rho(A_1, B_1)$	$\rho(A_2, B_1)$	$\rho(A_3, B_1)$	$\rho(A_4, B_1)$
B_2	$\rho(A_1, B_2)$	$\rho(A_2, B_2)$	$\rho(A_3, B_2)$	$\rho(A_4, B_2)$
\vdots	\vdots	\vdots	\vdots	\vdots
B_{24}	$\rho(A_1, B_{24})$	$\rho(A_2, B_{24})$	$\rho(A_3, B_{24})$	$\rho(A_4, B_{24})$

Benchmarks only focus on the LLMs’ performances in the specific skill they measure, like mathematical or coding ability. These skills, measured by benchmarks, relate to distinct task-specific rankings in LMArena. So we focus on the correlation coefficients between benchmark rankings and rankings in corresponding task-specific categories in LMArena. In detail, we pay more attention to the similarity between LLM rankings on benchmarks in the QA category with the overall ranking in LMArena, on benchmarks in the Mathematics category with the LMArena Math category, on benchmarks in the Code category with the LMArena Coding category, and on benchmarks in the Alignment category with the LMArena Instruction Following category.

4 RESULTS

This section presents the results of our research. Section 4.1 analyzes the correlation between benchmark rankings and rankings in corresponding task-specific categories in LMArena. We observe that some popular benchmarks have rankings not similar to LMArena’s. To further explore this, we investigate how the correlation coefficient relates to other factors—specifically, difficulty levels of benchmarks in Section 4.2 and release dates of benchmarks in Section 4.3. Finally, for each task-specific category in LMArena, Section 4.4 computes the average correlation coefficients across all associated benchmarks, concluding the section.

4.1 CORRELATION ANALYSIS OF INDIVIDUAL BENCHMARKS

The four tables Tables 3a to 3d present the correlation between each benchmark ranking and the ranking on the corresponding task-specific category in LMArena.

Table 3a shows that LLM ranking on FACTS Grounding shows no significant correlation with the overall ranking in LMArena, marked as “N.S.”. It also shows the lowest correlation ($\rho = 0.32$). Rankings on Humanity’s Last Exam also show no significant correlation, ranking second lowest in Table 3a. Table 3b shows that rankings on all mathematics benchmarks have a significant correlation with LMArena’s Math ranking of LLMs. Table 3c shows that ranking on IOI shows the lowest correlation ($\rho = 0.62$) with LMArena’s Code ranking of LLMs. Table 3d shows that ranking on IFEval shows a weak correlation ($\rho = 0.45$) with LMArena’s Instruction Following ranking of LLMs. Ranking on IFBench also shows a low correlation ($\rho = 0.61$, ranking second lowest in Table 3d). We noticed that both IFEval and IFBench use predefined rules for evaluation. Conversely, Creative Writing v3, Arena-Hard, and WritingBench use LLMs as judges. This suggests that Alignment benchmarks that use LLMs as judges may be better at resembling human perception.

Conclusion: Some popular benchmarks have rankings not similar to LMArena’s, like Humanity’s Last Exam, FACTS Grounding, and IFEval. LLM performance on these benchmarks may differ from human perception.

Table 3: The Spearman correlation coefficient (ρ) between each benchmark ranking and the ranking on the corresponding task-specific category in LMArena is presented in the table below. The table is organized into four sub-tables according to the type of benchmarks and their corresponding categories. Specifically, sub-table (a) shows the correlations for QA benchmarks against the overall ranking in LMArena; sub-table (b) for Mathematics benchmarks against the Math category in LMArena; sub-table (c) for Code benchmarks against the Coding category in LMArena; and sub-table (d) for Alignment benchmarks against the Instruction Following category in LMArena. N represents the number of LLMs that are common to both rankings being compared when calculating the Spearman correlation. For statistical significance, we use the p-value with thresholds: $***p < 0.001$, $**p < 0.01$, $*p < 0.05$; N.S. = Not significant ($p \geq 0.05$).

(a) QA vs LMArena’s overall ranking				(b) Mathematics vs LMArena’s Math ranking			
Benchmark	ρ	N	Significance	Benchmark	ρ	N	Significance
SuperGPQA	0.92	27	***	MATH-500	0.89	23	***
MMLU-Pro	0.86	42	***	FrontierMath	0.87	29	***
SimpleQA	0.79	23	***	HMMT February 2025	0.85	18	***
GPQA	0.79	28	***	AIME	0.78	29	***
ARC-AGI-2	0.71	18	**	MGSM	0.71	27	***
Humanity’s Last Exam	0.51	12	N.S.				
FACTS Grounding	0.32	18	N.S.				

(c) Code vs LMArena’s Code ranking				(d) Alignment vs LMArena’s Instruction Following ranking			
Benchmark	ρ	N	Significance	Benchmark	ρ	N	Significance
Aider polyglot	0.87	25	***	Creative Writing v3	0.84	26	***
SWE-bench Verified	0.8	13	**	Arena-Hard	0.73	18	***
Terminal-Bench	0.8	49	***	WritingBench	0.72	25	***
HumanEval	0.78	59	***	IFBench	0.61	47	***
LiveCodeBench	0.77	14	**	IFEval	0.45	20	*
SciCode	0.68	69	***				
IOI	0.64	10	*				

The conclusion above raises a critical question: Why do certain benchmark rankings—such as Humanity’s Last Exam—show such low correlation with human perception? A hypothesis is that these benchmarks often test very difficult skills. For example, Humanity’s Last Exam poses extremely complex problems aimed at AGI-level reasoning. Given the observed low correlation for notoriously difficult benchmarks, a natural question arises: What is the relationship between benchmarks’ difficulty and correlation coefficient with the corresponding LMArena category? The following section explores the problem.

4.2 RELATIONSHIP BETWEEN BENCHMARK DIFFICULTY AND CORRELATION WITH HUMAN PERCEPTION

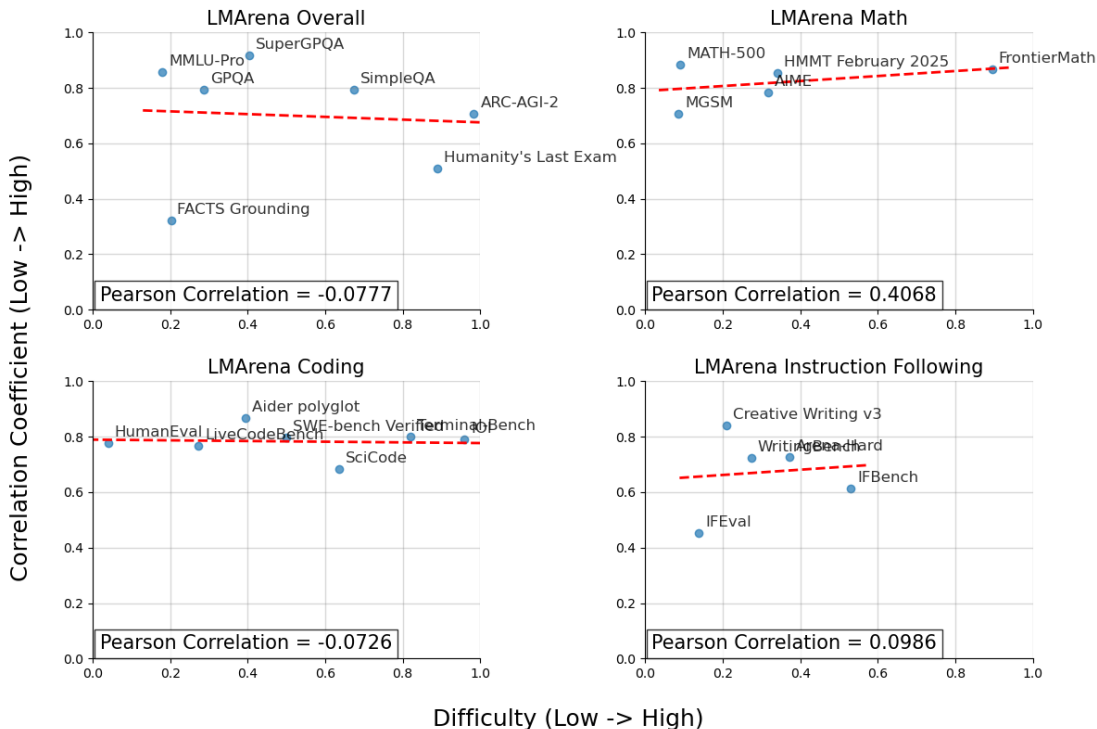


Figure 2: Relationship between benchmark difficulty and the Spearman correlation coefficient (ρ) with the corresponding LMArena category. The horizontal axis represents the average error rate of LLMs in a fixed set on the benchmark, where a higher value indicates greater difficulty. Each point (x_i, y_i) corresponds to a benchmark with an average error rate x_i , and y_i is the correlation coefficient between the LLM ranking on that benchmark and the ranking of the LMArena category shown above the subfigure. The red line is the regression line, and the indicated number is the Pearson correlation coefficient between x and y .

As discussed in the Evaluation sections of *Humanity's Last Exam* (Phan et al., 2025) and *Auto-Bench* (Li et al., 2025c), the accuracy achieved by top-tier LLMs on a benchmark can serve as an indicator of the benchmark's difficulty. To assess this difficulty more rigorously, we compute the average error rate, defined as one minus the average accuracy of a fixed set of LLMs. In our study, this set comprises models whose overall scores on the LMArena leaderboard fall between 1386 and 1418. We thus define benchmark difficulty as the average error rate. A higher difficulty value indicates a more challenging benchmark, which is likely to remain useful for evaluating future LLMs. Conversely, a lower difficulty suggests that the benchmark is nearing saturation and may have limited utility in discriminating among top-tier LLMs.

Figure 2 shows that benchmark difficulty does not strongly influence the Spearman rank-order correlation coefficient. This observation is consistent across all four categories of benchmarks, challenging a common assumption that more challenging benchmarks better reflect real-world performance. Our results do not support this view. This highlights the need for careful benchmark design that considers human preferences directly.

Conclusion: LLM performance on neither more difficult nor easier benchmarks is necessarily more similar to human perception.

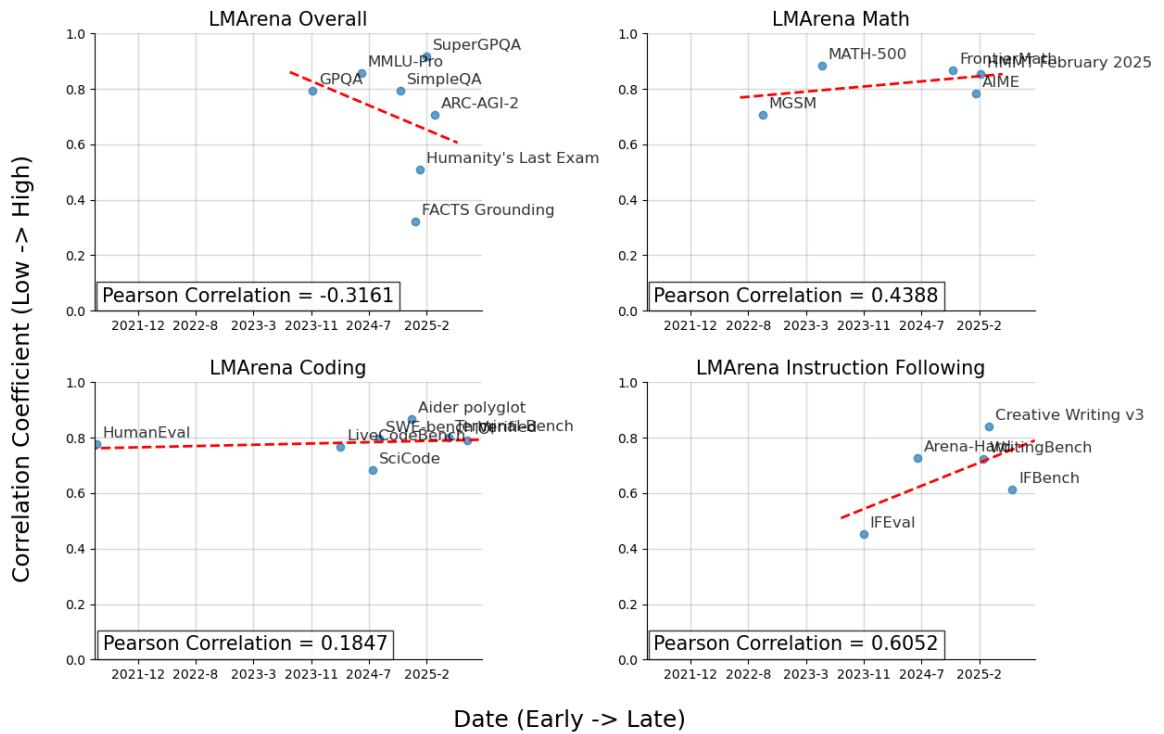


Figure 3: Relationship between benchmark publication date and the Spearman correlation coefficient (ρ) with the corresponding LMarena category. Each point (x_i, y_i) represents the benchmark that was released at date x_i , and y_i is the correlation coefficient between its ranking and the ranking of the LMarena category shown above the subfigure. The red line is the regression line, and the indicated number is the Pearson correlation coefficient between x and y .

4.3 RELATIONSHIP BETWEEN BENCHMARK RECENCY AND CORRELATION WITH HUMAN PERCEPTION

Figure 3 examines the relationship between the correlation coefficients of all benchmarks with the rankings on corresponding LMarena categories, and the initial release dates of these benchmarks. We create four scatter plots, with the horizontal axis representing the initial release date of benchmarks and the vertical axis showing the correlation coefficients between benchmarks and the corresponding LMarena categories.

For benchmarks within QA, Mathematics, and Code categories, the figures suggest a weak correlation between benchmarks' release date and their correlation with the corresponding LMarena category. For example, in the LMarena Overall subfigure, GPQA (released Nov 2023; $\rho = 0.79$) surpasses many subsequent benchmarks. In contrast, FACTS Grounding (released Jan 2025; $\rho = 0.32$) is inferior to those of earlier benchmarks. Conversely, in LMarena's Instruction Following subfigure, a positive trend between the correlation coefficients and the initial release dates is visible. Newer benchmarks like Creative Writing v3 (released Feb 2025; $\rho = 0.84$) show substantially higher correlation than older ones like IFEval (released Nov 2023; $\rho = 0.45$). These findings highlight a critical gap: While Alignment benchmarks are being improved to better resemble human perception of LLM outputs, QA, Mathematics, and Code benchmarks show no such progress. We call for future benchmarks within these categories to be more aligned with human perception.

Conclusion: LLM performance on newer benchmarks is not necessarily more similar to human perception, except for benchmarks under the Alignment category. Alignment benchmarks are being improved to better resemble human perception of LLM outputs.

Table 4: Average Spearman correlation coefficients (ρ) between ranking on the full set of 24 benchmarks and ranking on each of the four categories in LMArena.

	LMArena’s Overall Ranking	LMArena’s Math Ranking	LMArena’s Coding Ranking	LMArena’s Instruction Following Ranking
Average	0.69	0.78	0.69	0.69

4.4 AVERAGE CORRELATION COEFFICIENTS ACROSS LMARENA CATEGORIES

For each task-specific category in LMArena, Table 4 computes the average correlation coefficients across all associated benchmarks. This table is based on the schematic shown in Table 2, presenting the average of each column’s ρ . We find that the average calculated using the Math category in LMArena ($\rho = 0.78$) is the highest. This indicates that existing benchmarks generally can better evaluate human perception when using LLMs for solving math problems.

Conclusion: Compared with other abilities (e.g., coding ability and instruction following ability), the math ability of LLMs can be relatively better evaluated by current benchmarks.

5 DISCUSSION

In the previous analysis, we found that LLM performance on newer benchmarks within QA, Mathematics, and Code categories is not necessarily more similar to human perception. But we noticed τ^2 -Bench (Barres et al., 2025), which is a benchmark evaluating conversational agents in a dual-control environment, requiring coordination and communication between the agent and the user to succeed. This shows that capturing user perception has become an important target in the research community. But as this benchmark is the only one we found considering both reasoning ability and user perception, it’s not suitable in any of our 4 categories.

Looking ahead, new benchmarks in categories such as Question Answering and Code that more directly capture user perception can be designed. For example, a benchmark in the Code category should focus more on the readability, conciseness, and their ability to make basic modifications to user-written code while providing instructional guidance. By explicitly incorporating user perception into task design and evaluation protocols, the benchmark will provide a more faithful measure of LLM utility in practice and help bridge the gap between benchmark scores and real-world user perception.

6 CONCLUSION & LIMITATIONS

We systematically compared LLM performance on popular benchmarks with human perception as reflected in LMArena. The results indicate that several benchmarks, such as FACTS Grounding and IFEval, do not resemble human perception. Moreover, neither benchmark difficulty nor recency reliably predicts a benchmark’s resemblance to human perception. These findings highlight limitations in current benchmarking practices and the urgent need for more human-centered evaluation frameworks.

A key limitation of our study stems from LMArena itself. On LMArena, users often ask basic and simple questions, as more difficult or high-stakes problems are typically addressed directly using top-tier LLMs. Consequently, LMArena may not fully represent broader user perception. Additionally, some benchmarks include very few data points (e.g., IOI, with fewer than 15 LLMs being ranked on their leaderboards), which can introduce significant variability in correlation coefficients. For particularly difficult benchmarks where all LLMs achieve low accuracy, random factors may disproportionately affect performance, meaning small biases in accuracy can lead to large changes in LLM rankings and corresponding rank-order correlations. Finally, our categorization of benchmarks is necessarily approximate: most benchmarks evaluate multiple abilities. They cannot perfectly fit into a single broad category. For instance, Humanity’s Last Exam primarily tests multi-subject reasoning, yet some of its problems are mathematical in nature. These limitations suggest that our results, while indicative, should be interpreted with caution.

REFERENCES

- 486
487
488 Zhiqiang Shen Aidar Myrzakhan, Sondos Mahmoud Bsharat. Open-llm-leaderboard: From multi-
489 choice to open-style questions for llms evaluation, benchmark, and arena. *arXiv preprint*
490 *arXiv:2406.07545*, 2024.
- 491 Aider. o1 tops aider’s new polyglot leaderboard. [https://aider.chat/2024/12/](https://aider.chat/2024/12/21/polyglot.html#o1-tops-aider-s-new-polyglot-leaderboard)
492 [21/polyglot.html#o1-tops-aider-s-new-polyglot-leaderboard](https://aider.chat/2024/12/21/polyglot.html#o1-tops-aider-s-new-polyglot-leaderboard), 2024. Ac-
493 cessed: 2025-07-30.
- 494 Anthropic. Claude opus 4.1. <https://www.anthropic.com/news/claude-opus-4-1>,
495 2025. Accessed: 2025-09-21.
- 496 AoPS. 2025 aime i. [https://artofproblemsolving.com/wiki/index.php/2025_](https://artofproblemsolving.com/wiki/index.php/2025_AIME_I)
497 [AIME_I](https://artofproblemsolving.com/wiki/index.php/2025_AIME_I), 2025. Accessed: 2025-07-29.
- 499 Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. τ^2 -bench: Evaluating
500 conversational agents in a dual-control environment, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2506.07982)
501 [2506.07982](https://arxiv.org/abs/2506.07982).
- 502 Yixin Cao, Shibo Hong, Xinze Li, Jiahao Ying, Yubo Ma, Haiyuan Liang, Yantao Liu, Zijun Yao,
503 Xiaozhi Wang, Dan Huang, Wenxuan Zhang, Lifu Huang, Muhao Chen, Lei Hou, Qianru Sun,
504 Xingjun Ma, Zuxuan Wu, Min-Yen Kan, David Lo, Qi Zhang, Heng Ji, Jing Jiang, Juanzi Li,
505 Aixin Sun, Xuanjing Huang, Tat-Seng Chua, and Yu-Gang Jiang. Toward generalizable evaluation
506 in the llm era: A survey beyond benchmarks, 2025. URL [https://arxiv.org/abs/2504.](https://arxiv.org/abs/2504.18838)
507 [18838](https://arxiv.org/abs/2504.18838).
- 508 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared
509 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri,
510 Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,
511 Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,
512 Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fo-
513 tios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex
514 Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders,
515 Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec
516 Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob Mc-
517 Grew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large
518 language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- 519 Francois Chollet, Mike Knoop, Gregory Kamradt, Bryan Landers, and Henry Pinkard. Arc-agi-2:
520 A new challenge for frontier ai reasoning systems, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2505.11831)
521 [2505.11831](https://arxiv.org/abs/2505.11831).
- 522 Neil Chowdhury, James Aung, Chan Jun Shern, Oliver Jaffe, Dane Sherburn, Giulio Starace, Evan
523 Mays, Rachel Dias, Marwan Aljubei, Mia Glaese, Carlos E. Jimenez, John Yang, Leyton Ho,
524 Tejal Patwardhan, Kevin Liu, and Aleksander Madry. Introducing SWE-bench verified, 2024.
525 URL <https://openai.com/index/introducing-swe-bench-verified/>.
- 526 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
527 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin
528 Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-
529 Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric
530 Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornraphop Kawintiranon, Tania
531 Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilai Deutel, Nam Nguyen,
532 Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller,
533 Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero Rissa, Dan
534 Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves, Cindy
535 Wang, Nitesh Bharadwaj Gundavarapu, Iliia Shumailov, Bo Wang, Mantas Pajarskas, Joe Hey-
536 ward, Martin Nikolchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik,
537 Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kavukcuoglu,
538 Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov, Bo-
539 Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas Le-
ichner, Haichuan Yang, Zeld Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin

540 Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios
541 Vytiniotis, Jieru Mei, Mu Cai, Mohammed Badawi, Corey Fry, Ale Hartman, Daniel Zheng,
542 Eric Jia, James Keeling, Annie Louis, Ying Chen, Efen Robles, Wei-Chih Hung, Howard Zhou,
543 Nikita Saxena, Sonam Goenka, Olivia Ma, Zach Fisher, Mor Hazan Taege, Emily Graves, David
544 Steiner, Yujia Li, Sarah Nguyen, Rahul Sukthankar, Joe Stanton, Ali Eslami, Gloria Shen, Berkin
545 Akin, Alexey Guseynov, Yiqian Zhou, Jean-Baptiste Alayrac, Armand Joulin, Efrat Farkash,
546 Ashish Thapliyal, Stephen Roller, Noam Shazeer, Todor Davchev, Terry Koo, Hannah Forbes-
547 Pollard, Kartik Audhkhasi, Greg Farquhar, Adi Mayrav Gilady, Maggie Song, John Aslanides,
548 Piermaria Mendolicchio, Alicia Parrish, John Blitzer, Pramod Gupta, Xiaoen Ju, Xiaochen Yang,
549 Puranjay Datta, Andrea Tacchetti, Sanket Vaibhav Mehta, Gregory Dobb, Shubham Gupta, Feder-
550 erico Piccinini, Raia Hadsell, Sujee Rajayogam, Jiepu Jiang, Patrick Griffin, Patrik Sundberg,
551 Jamie Hayes, Alexey Frolov, Tian Xie, Adam Zhang, Kingshuk Dasgupta, Uday Kalra, Lior
552 Shani, Klaus Macherey, Tzu-Kuo Huang, Liam MacDermed, Karthik Duddu, Paulo Zacchello,
553 Zi Yang, Jessica Lo, Kai Hui, Matej Kastelic, Derek Gasaway, Qijun Tan, Summer Yue, Pablo
554 Barrio, John Wieting, Weel Yang, Andrew Nystrom, Solomon Demmessie, Anselm Levskaya,
555 Fabio Viola, Chetan Tekur, Greg Billock, George Necula, Mandar Joshi, Rylan Schaeffer, Swach-
556 hand Lokhande, Christina Sorokin, Pradeep Shenoy, Mia Chen, Mark Collier, Hongji Li, Tay-
557 lor Bos, Nevan Wichers, Sun Jae Lee, Angéline Pouget, Santhosh Thangaraj, Kyriakos Axio-
558 tis, Phil Crone, Rachel Sterneck, Nikolai Chinaev, Victoria Krakovna, Oleksandr Ferludin, Ian
559 Gemp, Stephanie Winkler, Dan Goldberg, Ivan Korotkov, Kefan Xiao, Malika Mehrotra, Sandeep
560 Mariserla, Vihari Piratla, Terry Thurk, Khiem Pham, Hongxu Ma, Alexandre Senges, Ravi Ku-
561 mar, Clemens Meyer, Ellie Talius, Nuo Wang Pierse, Ballie Sandhu, Horia Toma, Kuo Lin, Swa-
562 roop Nath, Tom Stone, Dorsa Sadigh, Nikita Gupta, Arthur Guez, Avi Singh, Matt Thomas, Tom
563 Duerig, Yuan Gong, Richard Tanburn, Lydia Lihui Zhang, Phuong Dao, Mohamed Hammad,
564 Sirui Xie, Shruti Rijhwani, Ben Murdoch, Duhyeon Kim, Will Thompson, Heng-Tze Cheng,
565 Daniel Sohn, Pablo Sprechmann, Qiantong Xu, Srinivas Tadepalli, Peter Young, Ye Zhang, Hansa
566 Srinivasan, Miranda Aperghis, Aditya Ayyar, Hen Fitoussi, Ryan Burnell, David Madras, Mike
567 Dusenberry, Xi Xiong, Tayo Oguntebi, Ben Albrecht, Jörg Bornschein, Jovana Mitrović, Ma-
568 son Dimarco, Bhargav Kanagal Shamanna, Premal Shah, Eren Sezener, Shyam Upadhyay, Dave
569 Lacey, Craig Schiff, Sebastian Baur, Sanjay Ganapathy, Eva Schnider, Mateo Wirth, Connor
570 Schenck, Andrey Simanovsky, Yi-Xuan Tan, Philipp Fränken, Dennis Duan, Bharath Mankalale,
571 Nikhil Dhawan, Kevin Sequeira, Zichuan Wei, Shivanker Goel, Caglar Unlu, Yukun Zhu, Haitian
572 Sun, Ananth Balashankar, Kurt Shuster, Megh Umekar, Mahmoud Alnahlawi, Aäron van den
573 Oord, Kelly Chen, Yuexiang Zhai, Zihang Dai, Kuang-Huei Lee, Eric Doi, Lukas Zilka, Rohith
574 Vallu, Disha Shrivastava, Jason Lee, Hisham Husain, Honglei Zhuang, Vincent Cohen-Addad,
575 Jarred Barber, James Atwood, Adam Sadovsky, Quentin Wellens, Steven Hand, Arunkumar Ra-
576 jendran, Aybuke Turker, CJ Carey, Yuanzhong Xu, Hagen Soltau, Zefei Li, Xinying Song, Cong-
577 long Li, Iurii Kemaev, Sasha Brown, Andrea Burns, Viorica Patraucean, Piotr Stanczyk, Renga
578 Aravamudhan, Mathieu Blondel, Hila Noga, Lorenzo Blanco, Will Song, Michael Isard, Mandar
579 Sharma, Reid Hayes, Dalia El Badawy, Avery Lamp, Itay Laish, Olga Kozlova, Kelvin Chan,
580 Sahil Singla, Srinivas Sunkara, Mayank Upadhyay, Chang Liu, Aijun Bai, Jarek Wilkiewicz,
581 Martin Zlocha, Jeremiah Liu, Zhuowan Li, Haiguang Li, Omer Barak, Ganna Raboshchuk, Jiho
582 Choi, Fangyu Liu, Erik Jue, Mohit Sharma, Andrea Marzoca, Robert Busa-Fekete, Anna Ko-
583 rsun, Andre Elisseeff, Zhe Shen, Sara Mc Carthy, Kay Lamerigts, Anahita Hosseini, Hanzhao
584 Lin, Charlie Chen, Fan Yang, Kushal Chauhan, Mark Omernick, Dawei Jia, Karina Zainul-
585 lina, Demis Hassabis, Danny Vainstein, Ehsan Amid, Xiang Zhou, Ronny Votel, Eszter Vértés,
586 Xinjian Li, Zongwei Zhou, Angeliki Lazaridou, Brendan McMahan, Arjun Narayanan, Hubert
587 Soyer, Sujoy Basu, Kayi Lee, Bryan Perozzi, Qin Cao, Leonard Berrada, Rahul Arya, Ke Chen,
588 Katrina, Xu, Matthias Lochbrunner, Alex Hofer, Sahand Sharifzadeh, Renjie Wu, Sally Gold-
589 man, Pranjal Awasthi, Xuezhi Wang, Yan Wu, Claire Sha, Biao Zhang, Maciej Mikula, Filippo
590 Graziano, Siobhan Mcloughlin, Irene Giannoumis, Youhei Namiki, Chase Malik, Carey Rade-
591 baugh, Jamie Hall, Ramiro Leal-Cavazos, Jianmin Chen, Vikas Sindhwani, David Kao, David
592 Greene, Jordan Griffith, Chris Welty, Ceslee Montgomery, Toshihiro Yoshino, Liangzhe Yang,
593 Noah Goodman, Assaf Hurwitz Michaely, Kevin Lee, KP Sawhney, Wei Chen, Zheng Zheng,
Megan Shum, Nikolay Savinov, Etienne Pot, Alex Pak, Morteza Zadimoghaddam, Sijal Bhat-
nagar, Yoad Lewenberg, Blair Kutzman, Ji Liu, Lesley Katzen, Jeremy Selier, Josip Djolonga,
Dmitry Lepikhin, Kelvin Xu, Jacky Liang, Jiewen Tan, Benoit Schillings, Muge Ersoy, Pete
Blois, Bernd Bandemer, Abhimanyu Singh, Sergei Lebedev, Pankaj Joshi, Adam R. Brown,
Evan Palmer, Shreya Pathak, Komal Jalan, Fedir Zubach, Shuba Lall, Randall Parker, Alok Gun-

594 jan, Sergey Rogulenko, Sumit Sanghai, Zhaoqi Leng, Zoltan Egyed, Shixin Li, Maria Ivanova,
 595 Kostas Andriopoulos, Jin Xie, Elan Rosenfeld, Auriel Wright, Ankur Sharma, Xinyang Geng,
 596 Yicheng Wang, Sam Kwei, Renke Pan, Yujing Zhang, Gabby Wang, Xi Liu, Chak Yeung, Eliz-
 597 abeth Cole, Aviv Rosenberg, Zhen Yang, Phil Chen, George Polovets, Pranav Nair, Rohun Sax-
 598 ena, Josh Smith, Shuo yiin Chang, Aroma Mahendru, Svetlana Grant, Anand Iyer, Irene Cai,
 599 Jed McGiffin, Jiaming Shen, Alanna Walton, Antonious Girgis, Oliver Woodman, Rosemary Ke,
 600 Mike Kwong, Louis Rouillard, Jinneng Rao, Zhihao Li, Yuntao Xu, Flavien Prost, Chi Zou,
 601 Ziwei Ji, Alberto Magni, Tyler Liechty, Dan A. Calian, Deepak Ramachandran, Igor Krivokon,
 602 Hui Huang, Terry Chen, Anja Hauth, Anastasija Ilić, Weijuan Xi, Hyeontaek Lim, Vlad-Doru
 603 Ion, Pooya Moradi, Metin Toksoz-Exley, Kalesha Bullard, Miltos Allamanis, Xiaomeng Yang,
 604 Sophie Wang, Zhi Hong, Anita Gergely, Cheng Li, Bhavishya Mittal, Vitaly Kovalev, Victor Un-
 605 gureanu, Jane Labanowski, Jan Wassenberg, Nicolas Lacasse, Geoffrey Cideron, Petar Dević,
 606 Annie Marsden, Lynn Nguyen, Michael Fink, Yin Zhong, Tatsuya Kiyono, Desi Ivanov, Sally
 607 Ma, Max Bain, Kiran Yalasang, Jennifer She, Anastasia Petrushkina, Mayank Lunayach, Carla
 608 Bromberg, Sarah Hodgkinson, Vilobh Meshram, Daniel Vlasic, Austin Kyker, Steve Xu, Jeff Stan-
 609 way, Zuguang Yang, Kai Zhao, Matthew Tung, Seth Odoom, Yasuhisa Fujii, Justin Gilmer,
 610 Eunyoung Kim, Felix Halim, Quoc Le, Bernd Bohnet, Seliem El-Sayed, Behnam Neyshabur,
 611 Malcolm Reynolds, Dean Reich, Yang Xu, Erica Moreira, Anuj Sharma, Zeyu Liu, Moham-
 612 mad Javad Hosseini, Naina Raisinghani, Yi Su, Ni Lao, Daniel Formoso, Marco Gelmi, Almog
 613 Gueta, Tapomay Dey, Elena Gribovskaya, Domagoj Čevič, Sidharth Mudgal, Garrett Bingham,
 614 Jianling Wang, Anurag Kumar, Alex Cullum, Feng Han, Konstantinos Bousmalis, Diego Cedillo,
 615 Grace Chu, Vladimir Magay, Paul Michel, Ester Hlavnova, Daniele Calandriello, Setareh Ari-
 616 afar, Kaisheng Yao, Vikash Schwag, Arpi Vezer, Agustin Dal Lago, Zhenkai Zhu, Paul Kishan
 617 Rubenstein, Allen Porter, Anirudh Baddepudi, Oriana Riva, Mihai Dorin Istin, Chih-Kuan Yeh,
 618 Zhi Li, Andrew Howard, Nilpa Jha, Jeremy Chen, Raoul de Liedekerke, Zafarali Ahmed, Mikel
 619 Rodriguez, Tanuj Bhatia, Bangju Wang, Ali Elqursh, David Klinghoffer, Peter Chen, Pushmeet
 620 Kohli, Te I, Weiyang Zhang, Zack Nado, Jilin Chen, Maxwell Chen, George Zhang, Aayush
 621 Singh, Adam Hillier, Federico Lebron, Yiqing Tao, Ting Liu, Gabriel Dulac-Arnold, Jingwei
 622 Zhang, Shashi Narayan, Bu Huang Liu, Orhan Firat, Abhishek Bhowmick, Bingyuan Liu, Hao
 623 Zhang, Zizhao Zhang, Georges Rotival, Nathan Howard, Anu Sinha, Alexander Grushetsky, Ben-
 624 jamin Beyret, Keerthana Gopalakrishnan, James Zhao, Kyle He, Szabolcs Payrits, Zaid Nabulsi,
 625 Zhaoyi Zhang, Weijie Chen, Edward Lee, Nova Fallen, Sreenivas Gollapudi, Aurick Zhou, Filip
 626 Pavetić, Thomas Köppe, Shiyu Huang, Rama Pasumarthi, Nick Fernando, Felix Fischer, Daria
 627 Ćurko, Yang Gao, James Svensson, Austin Stone, Haroon Qureshi, Abhishek Sinha, Apoorv
 628 Kulshreshtha, Martin Matysiak, Jieming Mao, Carl Saroufim, Aleksandra Faust, Qingnan Duan,
 629 Gil Fidel, Kaan Katircioglu, Raphaël Lopez Kaufman, Dhruv Shah, Weize Kong, Abhishek
 630 Bapna, Gellért Weisz, Emma Dunleavy, Praneet Dutta, Tianqi Liu, Rahma Chaabouni, Carolina
 631 Parada, Marcus Wu, Alexandra Belias, Alessandro Bissacco, Stanislav Fort, Li Xiao, Fantine
 632 Huot, Chris Knutsen, Yochai Blau, Gang Li, Jennifer Prendki, Juliette Love, Yinlam Chow, Pichi
 633 Charoenpanit, Hidetoshi Shimokawa, Vincent Coriou, Karol Gregor, Tomas Izo, Arjun Akula,
 634 Mario Pinto, Chris Hahn, Dominik Paulus, Jiaxian Guo, Neha Sharma, Cho-Jui Hsieh, Adaeze
 635 Chukwuka, Kazuma Hashimoto, Nathalie Rauschmayr, Ling Wu, Christof Angermueller, Yulong
 636 Wang, Sebastian Gerlach, Michael Pliskin, Daniil Mirylenka, Min Ma, Lexi Baugher, Bryan Gale,
 637 Shaan Bijwadia, Nemanja Rakićević, David Wood, Jane Park, Chung-Ching Chang, Babi Seal,
 638 Chris Tar, Kacper Krasowiak, Yiwen Song, Georgi Stephanov, Gary Wang, Marcello Maggioni,
 639 Stein Xudong Lin, Felix Wu, Shachi Paul, Zixuan Jiang, Shubham Agrawal, Bilal Piot, Alex
 640 Feng, Cheolmin Kim, Tulsee Doshi, Jonathan Lai, Chuqiao, Xu, Sharad Vikram, Ciprian Chelba,
 641 Sebastian Krause, Vincent Zhuang, Jack Rae, Timo Denk, Adrian Collister, Lotte Weerts, Xi-
 642 anghong Luo, Yifeng Lu, Håvard Garnes, Nitish Gupta, Terry Spitz, Avinatan Hassidim, Lihao
 643 Liang, Izhak Shafran, Peter Humphreys, Kenny Vassigh, Phil Wallis, Virat Shejwalkar, Nicolas
 644 Perez-Nieves, Rachel Hornung, Melissa Tan, Beka Westberg, Andy Ly, Richard Zhang, Brian
 645 Farris, Jongbin Park, Alec Kosik, Zeynep Cankara, Andrii Maksai, Yunhan Xu, Albin Cassirer,
 646 Sergi Caelles, Abbas Abdolmaleki, Mencher Chiang, Alex Fabrikant, Shravya Shetty, Luheng
 647 He, Mai Giménez, Hadi Hashemi, Sheena Panthaplackel, Yana Kulizhskaya, Salil Deshmukh,
 Daniele Pighin, Robin Alazard, Disha Jindal, Seb Noury, Pradeep Kumar S, Siyang Qin, Xerxes
 Dotiwalla, Stephen Spencer, Mohammad Babaeizadeh, Blake Jianhang Chen, Vaibhav Mehta,
 Jennie Lees, Andrew Leach, Penporn Koanantakool, Ilia Akolzin, Ramona Comanescu, Junwhan
 Ahn, Alexey Svyatkovskiy, Basil Mustafa, David D’Ambrosio, Shiva Mohan Reddy Garlapati,

648 Pascal Lamblin, Alekh Agarwal, Shuang Song, Pier Giuseppe Sessa, Pauline Coquinot, John
649 Maggs, Hussain Masoom, Divya Pitta, Yaqing Wang, Patrick Morris-Suzuki, Billy Porter, John-
650 son Jia, Jeffrey Dudek, Raghavender R, Cosmin Paduraru, Alan Ansell, Tolga Bolukbasi, Tony
651 Lu, Ramya Ganeshan, Zi Wang, Henry Griffiths, Rodrigo Benenson, Yifan He, James Swirhun,
652 George Papamakarios, Aditya Chawla, Kuntal Sengupta, Yan Wang, Vedrana Milutinovic, Igor
653 Mordatch, Zhipeng Jia, Jamie Smith, Will Ng, Shitij Nigam, Matt Young, Eugen Vušak, Blake
654 Hechtman, Sheela Goenka, Avital Zipori, Kareem Ayoub, Ashok Papat, Trilok Acharya, Luo
655 Yu, Dawn Bloxwich, Hugo Song, Paul Roit, Haiqiong Li, Aviel Boag, Nigamaa Nayakanti,
656 Bilva Chandra, Tianli Ding, Aahil Mehta, Cath Hope, Jiageng Zhang, Idan Heimlich Shtacher,
657 Kartikeya Badola, Ryo Nakashima, Andrei Sozanschi, Iulia Comşa, Ante Žužul, Emily Cave-
658 ness, Julian Odell, Matthew Watson, Dario de Cesare, Phillip Lippe, Derek Lockhart, Siddharth
659 Verma, Huizhong Chen, Sean Sun, Lin Zhuo, Aditya Shah, Prakhar Gupta, Alex Muzio, Ning
660 Niu, Amir Zait, Abhinav Singh, Meenu Gaba, Fan Ye, Prajit Ramachandran, Mohammad Saleh,
661 Raluca Ada Popa, Ayush Dubey, Frederick Liu, Sara Javanmardi, Mark Epstein, Ross Hemsley,
662 Richard Green, Nishant Ranka, Eden Cohen, Chuyuan Kelly Fu, Sanjay Ghemawat, Jed Borovik,
663 James Martens, Anthony Chen, Pranav Shyam, André Susano Pinto, Ming-Hsuan Yang, Alexan-
664 dru Țifrea, David Du, Boqing Gong, Ayushi Agarwal, Seungyeon Kim, Christian Frank, Saloni
665 Shah, Xiaodan Song, Zhiwei Deng, Ales Mikhalap, Kleopatra Chatziprimou, Timothy Chung,
666 Toni Creswell, Susan Zhang, Yennie Jun, Carl Lebsack, Will Truong, Slavica Andaćić, Itay Yona,
667 Marco Fornoni, Rong Rong, Serge Topov, Afzal Shama Soudagar, Andrew Audibert, Salah
668 Zaiem, Zaheer Abbas, Andrei Rusu, Sahitya Potluri, Shitao Weng, Anastasios Kementsietsidis,
669 Anton Tsitsulin, Daiyi Peng, Natalie Ha, Sanil Jain, Tejasi Latkar, Simeon Ivanov, Cory McLean,
670 Anirudh GP, Rajesh Venkataraman, Canoe Liu, Dilip Krishnan, Joel D’sa, Roey Yogev, Paul
671 Collins, Benjamin Lee, Lewis Ho, Carl Doersch, Gal Yona, Shawn Gao, Felipe Tiengo Ferreira,
672 Adnan Ozturk, Hannah Muckenhirn, Ce Zheng, Gargi Balasubramaniam, Mudit Bansal, George
673 van den Driessche, Sivan Eiger, Salem Haykal, Vedant Misra, Abhimanyu Goyal, Danilo Martins,
674 Gary Leung, Jonas Valfridsson, Four Flynn, Will Bishop, Chenxi Pang, Yoni Halpern, Honglin
675 Yu, Lawrence Moore, Yuvein, Zhu, Sridhar Thiagarajan, Yoel Drori, Zhisheng Xiao, Lucio Dery,
676 Rolf Jagerman, Jing Lu, Eric Ge, Vaibhav Aggarwal, Arjun Khare, Vinh Tran, Oded Elyada,
677 Ferran Alet, James Rubin, Ian Chou, David Tian, Libin Bai, Lawrence Chan, Lukas Lew, Karolis
678 Misiunas, Taylan Bilal, Aniket Ray, Sindhu Raghuram, Alex Castro-Ros, Viral Carpenter,
679 CJ Zheng, Michael Kilgore, Josef Broder, Emily Xue, Praveen Kallakuri, Dheeru Dua, Nancy
680 Yuen, Steve Chien, John Schultz, Saurabh Agrawal, Reut Tsarfaty, Jingcao Hu, Ajay Kannan,
681 Dror Marcus, Nisarg Kothari, Baochen Sun, Ben Horn, Matko Bošnjak, Ferjad Naem, Dean
682 Hirsch, Lewis Chiang, Boya Fang, Jie Han, Qifei Wang, Ben Hora, Antoine He, Mario Lučić,
683 Beer Changpinyo, Anshuman Tripathi, John Youssef, Chester Kwak, Philippe Schlattner, Cat
684 Graves, Rémi Leblond, Wenjun Zeng, Anders Andreassen, Gabriel Rasskin, Yue Song, Eddie
685 Cao, Junhyuk Oh, Matt Hoffman, Wojtek Skut, Yichi Zhang, Jon Stritar, Xingyu Cai, Saarthak
686 Khanna, Kathie Wang, Shriya Sharma, Christian Reisswig, Younghoon Jun, Aman Prasad, Ta-
687 tiana Sholokhova, Preeti Singh, Adi Gerzi Rosenthal, Anian Ruoss, Françoise Beaufays, Sean
688 Kirmani, Dongkai Chen, Johan Schalkwyk, Jonathan Herzig, Been Kim, Josh Jacob, Damien
689 Vincent, Adrian N Reyes, Ivana Balazevic, Léonard Hussenot, Jon Schneider, Parker Barnes,
690 Luis Castro, Spandana Raj Babbula, Simon Green, Serkan Cabi, Nico Duduta, Danny Driess,
691 Rich Galt, Noam Velan, Junjie Wang, Hongyang Jiao, Matthew Mauger, Du Phan, Miteyan Patel,
692 Vlado Galić, Jerry Chang, Eyal Marcus, Matt Harvey, Julian Salazar, Elahe Dabir, Suraj Satishku-
693 mar Sheth, Amol Mandhane, Hanie Sedghi, Jeremiah Willcock, Amir Zandieh, Shruthi Prab-
694 hakara, Aida Amini, Antoine Miech, Victor Stone, Massimo Nicosia, Paul Niemczyk, Ying Xiao,
695 Lucy Kim, Sławek Kwasiborski, Vikas Verma, Ada Maksutaj Oflazer, Christoph Hirsenschall, Pe-
696 ter Sung, Lu Liu, Richard Everett, Michiel Bakker, Ágoston Weisz, Yufei Wang, Vivek Sam-
697 pathkumar, Uri Shaham, Bibo Xu, Yasemin Altun, Mingqiu Wang, Takaaki Saeki, Guanjie Chen,
698 Emanuel Taropa, Shanthal Vasanth, Sophia Austin, Lu Huang, Goran Petrovic, Qingyun Dou,
699 Daniel Golovin, Grigory Rozhdestvenskiy, Allie Culp, Will Wu, Motoki Sano, Divya Jain, Julia
700 Proskurnia, Sébastien Cevey, Alejandro Cruzado Ruiz, Piyush Patil, Mahdi Mirzazadeh, Eric Ni,
701 Javier Snaider, Lijie Fan, Alexandre Fréchet, AJ Pierigiovanni, Shariq Iqbal, Kenton Lee, Clau-
dio Fantacci, Jinwei Xing, Lisa Wang, Alex Irgan, David Raposo, Yi Luan, Zhuoyuan Chen, Har-
ish Ganapathy, Kevin Hui, Jiazhong Nie, Isabelle Guyon, Heming Ge, Roopali Vij, Hui Zheng,
Dayeong Lee, Alfonso Castaño, Khuslen Baatarsukh, Gabriel Ibagon, Alexandra Chronopoulou,
Nicholas FitzGerald, Shashank Viswanadha, Safeen Huda, Rivka Moroshko, Georgi Stoyanov,
Prateek Kolhar, Alain Vaucher, Ishaan Watts, Adhi Kuncoro, Henryk Michalewski, Satish Kam-

702 bala, Bat-Orgil Batsaikhan, Alek Andreev, Irina Jurenka, Maigo Le, Qihang Chen, Wael Al Jishi,
703 Sarah Chakera, Zhe Chen, Aditya Kini, Vikas Yadav, Aditya Siddhant, Iliia Labzovsky, Balaji
704 Lakshminarayanan, Carrie Grimes Bostock, Pankil Botadra, Ankesh Anand, Colton Bishop, Sam
705 Conway-Rahman, Mohit Agarwal, Yani Donchev, Achintya Singhal, Félix de Chaumont Quiry,
706 Natalia Ponomareva, Nishant Agrawal, Bin Ni, Kalpesh Krishna, Masha Samsikova, John Karro,
707 Yilun Du, Tamara von Glehn, Caden Lu, Christopher A. Choquette-Choo, Zhen Qin, Tingnan
708 Zhang, Sicheng Li, Divya Tyam, Swaroop Mishra, Wing Lowe, Colin Ji, Weiyi Wang, Man-
709 aal Faruqi, Ambrose Slone, Valentin Dalibard, Arunachalam Narayanaswamy, John Lambert,
710 Pierre-Antoine Manzagol, Dan Karliner, Andrew Bolt, Ivan Lobov, Aditya Kusupati, Chang
711 Ye, Xuan Yang, Heiga Zen, Nelson George, Mukul Bhutani, Olivier Lacombe, Robert Riachi,
712 Gagan Bansal, Rachel Soh, Yue Gao, Yang Yu, Adams Yu, Emily Nottage, Tania Rojas-Esponda,
713 James Noraky, Manish Gupta, Ragha Kotikalapudi, Jichuan Chang, Sanja Deur, Dan Graur, Alex
714 Mossin, Erin Farnese, Ricardo Figueira, Alexandre Moufarek, Austin Huang, Patrik Zochbauer,
715 Ben Ingram, Tongzhou Chen, Zelin Wu, Adrià Puigdomènech, Leland Rechis, Da Yu, Sri Gay-
716 atri Sundara Padmanabhan, Rui Zhu, Chu ling Ko, Andrea Banino, Samira Daruki, Aarush Sel-
717 van, Dhruva Bhaswar, Daniel Hernandez Diaz, Chen Su, Salvatore Scellato, Jennifer Brennan,
718 Woohyun Han, Grace Chung, Priyanka Agrawal, Urvashi Khandelwal, Khe Chai Sim, Mor-
719 gane Lustman, Sam Ritter, Kelvin Guu, Jiawei Xia, Prateek Jain, Emma Wang, Tyrone Hill,
720 Mirko Rossini, Marija Kostelac, Tautvydas Misiunas, Amit Sabne, Kyuyeun Kim, Ahmet Is-
721 cen, Congchao Wang, José Leal, Ashwin Sreevatsa, Utku Evci, Manfred Warmuth, Saket Joshi,
722 Daniel Suo, James Lottes, Garrett Honke, Brendan Jou, Stefani Karp, Jieru Hu, Himanshu Sahni,
723 Adrien Ali Taïga, William Kong, Samrat Ghosh, Renshen Wang, Jay Pavagadhi, Natalie Axels-
724 son, Nikolai Grigorev, Patrick Siegler, Rebecca Lin, Guohui Wang, Emilio Parisotto, Sharath
725 Maddineni, Krishan Subudhi, Eyal Ben-David, Elena Pochernina, Orgad Keller, Thi Avrahami,
726 Zhe Yuan, Pulkit Mehta, Jialu Liu, Sherry Yang, Wendy Kan, Katherine Lee, Tom Funkhouser,
727 Derek Cheng, Hongzhi Shi, Archit Sharma, Joe Kelley, Matan Eyal, Yury Malkov, Corentin Tal-
728 lec, Yuval Bahat, Shen Yan, Xintian, Wu, David Lindner, Chengda Wu, Avi Caciularu, Xiyang
729 Luo, Rodolphe Jenatton, Tim Zaman, Yingying Bi, Ilya Kornakov, Ganesh Mallya, Daisuke
730 Ikeda, Itay Karo, Anima Singh, Colin Evans, Praneeth Netrapalli, Vincent Nallatamby, Isaac
731 Tian, Yanniss Assael, Vikas Raunak, Victor Carbune, Ioana Bica, Lior Madmoni, Dee Cattle,
732 Snchit Grover, Krishna Somandepalli, Sid Lall, Amelio Vázquez-Reina, Riccardo Patana, Jiaqi
733 Mu, Pranav Talluri, Maggie Tran, Rajeev Aggarwal, RJ Skerry-Ryan, Jun Xu, Mike Burrows,
734 Xiaoyue Pan, Edouard Yvinec, Di Lu, Zhiying Zhang, Duc Dung Nguyen, Hairong Mu, Gabriel
735 Barcik, Helen Ran, Lauren Beltrone, Krzysztof Choromanski, Dia Kharrat, Samuel Albanie, Sean
736 Purser-haskell, David Bieber, Carrie Zhang, Jing Wang, Tom Hudson, Zhiyuan Zhang, Han Fu,
737 Johannes Mauerer, Mohammad Hossein Bateni, AJ Maschinot, Bing Wang, Muye Zhu, Arjun
738 Pillai, Tobias Weyand, Shuang Liu, Oscar Akerlund, Fred Bertsch, Vittal Premachandran, Ali-
739 cia Jin, Vincent Roulet, Peter de Boursac, Shubham Mittal, Ndaba Ndebele, Georgi Karadzhov,
740 Sahra Ghalebikesabi, Ricky Liang, Allen Wu, Yale Cong, Nimesh Ghelani, Sumeet Singh, Ba-
741 har Fatemi, Warren, Chen, Charles Kwong, Alexey Kolganov, Steve Li, Richard Song, Chenkai
742 Kuang, Sobhan Miryoosefi, Dale Webster, James Wendt, Arkadiusz Socala, Guolong Su, Artur
743 Mendonça, Abhinav Gupta, Xiaowei Li, Tomy Tsai, Qiong, Hu, Kai Kang, Angie Chen, Ser-
744 tan Girgin, Yongqin Xian, Andrew Lee, Nolan Ramsden, Leslie Baker, Madeleine Clare Elish,
745 Varvara Krayvanova, Rishabh Joshi, Jiri Simsa, Yao-Yuan Yang, Piotr Ambroszczyk, Dipankar
746 Ghosh, Arjun Kar, Yuan Shangguan, Yumeya Yamamori, Yaroslav Akulov, Andy Brock, Hao-
747 tian Tang, Siddharth Vashishtha, Rich Munoz, Andreas Steiner, Kalyan Andra, Daniel Eppens,
748 Qixuan Feng, Hayato Kobayashi, Sasha Goldshtein, Mona El Mahdy, Xin Wang, Jilei, Wang,
749 Richard Killam, Tom Kwiatkowski, Kavaya Koppurapu, Serena Zhan, Chao Jia, Alexei Bende-
750 bury, Sheryl Luo, Adrià Recasens, Timothy Knight, Jing Chen, Mohak Patel, YaGuang Li, Ben
751 Withbroe, Dean Weesner, Kush Bhatia, Jie Ren, Danielle Eisenbud, Ebrahim Songhori, Yanhua
752 Sun, Travis Choma, Tasos Kementsietsidis, Lucas Manning, Brian Roark, Wael Farhan, Jie Feng,
753 Susheel Tatineni, James Cobon-Kerr, Yunjie Li, Lisa Anne Hendricks, Isaac Noble, Chris Breaux,
754 Nate Kushman, Liqian Peng, Fuzhao Xue, Taylor Tobin, Jamie Rogers, Josh Lipschultz, Chris
755 Alberti, Alexey Vlaskin, Mostafa Dehghani, Roshan Sharma, Tris Warkentin, Chen-Yu Lee, Be-
756 nigno Urias, Da-Cheng Juan, Angad Chandorkar, Hila Sheftel, Ruibo Liu, Elnaz Davoodi, Borja
757 De Balle Pigem, Kedar Dhamdhere, David Ross, Jonathan Hoech, Mahdis Mahdieh, Li Liu, Qiu-
758 jia Li, Liam McCafferty, Chenxi Liu, Markus Mircea, Yunting Song, Omkar Savant, Alaa Saade,
759 Colin Cherry, Vincent Hellendoorn, Siddharth Goyal, Paul Pucciarelli, David Vilar Torres, Zo-
760 har Yahav, Hyo Lee, Lars Lowe Sjoesund, Christo Kirov, Bo Chang, Deepanway Ghoshal, Lu Li,

756 Gilles Baechler, Sébastien Pereira, Tara Sainath, Anudhyan Boral, Dominik Grewe, Afief Halumi,
 757 Nguyet Minh Phu, Tianxiao Shen, Marco Tulio Ribeiro, Dhriti Varma, Alex Kaskasoli, Vlad Fein-
 758 berg, Navneet Potti, Jarrod Kahn, Matheus Wisniewski, Shakir Mohamed, Arnar Mar Hrafnkels-
 759 son, Bobak Shahriari, Jean-Baptiste Lespiau, Lisa Patel, Legg Yeung, Tom Paine, Lantao Mei,
 760 Alex Ramirez, Rakesh Shivanna, Li Zhong, Josh Woodward, Guilherme Tubone, Samira Khan,
 761 Heng Chen, Elizabeth Nielsen, Catalin Ionescu, Utsav Prabhu, Mingcen Gao, Qingze Wang, Sean
 762 Augenstein, Neesha Subramaniam, Jason Chang, Fotis Iliopoulos, Jiaming Luo, Myriam Khan,
 763 Weicheng Kuo, Denis Teplyashin, Florence Perot, Logan Kilpatrick, Amir Globerson, Hongkun
 764 Yu, Anfal Siddiqui, Nick Sukhanov, Arun Kandoor, Umang Gupta, Marco Andreetto, Moran
 765 Ambar, Donnie Kim, Paweł Wesolowski, Sarah Perrin, Ben Limonchik, Wei Fan, Jim Stephan,
 766 Ian Stewart-Binks, Ryan Kappedal, Tong He, Sarah Cogan, Romina Datta, Tong Zhou, Jiayu
 767 Ye, Leandro Kieliger, Ana Ramalho, Kyle Kastner, Fabian Mentzer, Wei-Jen Ko, Arun Suggala,
 768 Tianhao Zhou, Shiraz Butt, Hana Strejček, Lior Belenki, Subhashini Venugopalan, Mingyang
 769 Ling, Evgenii Eltyshv, Yunxiao Deng, Geza Kovacs, Mukund Raghavachari, Hanjun Dai, Tal
 770 Schuster, Steven Schwarcz, Richard Nguyen, Arthur Nguyen, Gavin Buttimore, Shrestha Basu
 771 Mallick, Sudeep Gandhe, Seth Benjamin, Michal Jastrzebski, Le Yan, Sugato Basu, Chris Apps,
 772 Isabel Edkins, James Allingham, Immanuel Odisho, Tomas Kocisky, Jewel Zhao, Linting Xue,
 773 Apoorv Reddy, Chrysovalantis Anastasiou, Aviel Atlas, Sam Redmond, Kieran Milan, Nicolas
 774 Heess, Herman Schmit, Allan Dafoe, Daniel Andor, Tynan Gangwani, Anca Dragan, Sheng
 775 Zhang, Ashyana Kachra, Gang Wu, Siyang Xue, Kevin Aydin, Siqi Liu, Yuxiang Zhou, Mahan
 776 Malihi, Austin Wu, Siddharth Gopal, Candice Schumann, Peter Stys, Alek Wang, Mirek Olšák,
 777 Dangyi Liu, Christian Schallhart, Yiran Mao, Demetra Brady, Hao Xu, Tomas Mery, Chawin
 778 Sitawarin, Siva Velusamy, Tom Cobley, Alex Zhai, Christian Walder, Nitzan Katz, Ganesh Jawa-
 779 har, Chinmay Kulkarni, Antoine Yang, Adam Paszke, Yinan Wang, Bogdan Damoc, Zalán Bor-
 780 sos, Ray Smith, Jinning Li, Mansi Gupta, Andrei Kapishnikov, Sushant Prakash, Florian Luisier,
 781 Rishabh Agarwal, Will Grathwohl, Kuangyuan Chen, Kehang Han, Nikhil Mehta, Andrew Over,
 782 Shekoofeh Azizi, Lei Meng, Niccolò Dal Santo, Kelvin Zheng, Jane Shapiro, Igor Petrovski,
 783 Jeffrey Hui, Amin Ghafouri, Jasper Snoek, James Qin, Mandy Jordan, Caitlin Sikora, Jonathan
 784 Malmaud, Yuheng Kuang, Aga Świetlik, Ruoxin Sang, Chongyang Shi, Leon Li, Andrew Rosen-
 785 berg, Shubin Zhao, Andy Crawford, Jan-Thorsten Peter, Yun Lei, Xavier Garcia, Long Le, Todd
 786 Wang, Julien Amelot, Dave Orr, Praneeth Kacham, Dana Alon, Gladys Tyen, Abhinav Arora,
 787 James Lyon, Alex Kurakin, Mimi Ly, Theo Guidroz, Zhipeng Yan, Rina Panigrahy, Pingmei
 788 Xu, Thais Kagohara, Yong Cheng, Eric Noland, Jinhyuk Lee, Jonathan Lee, Cathy Yip, Maria
 789 Wang, Efrat Nehoran, Alexander Bykovsky, Zhihao Shan, Ankit Bhagatwala, Chaochao Yan, Jie
 790 Tan, Guillermo Garrido, Dan Ethier, Nate Hurley, Grace Vesom, Xu Chen, Siyuan Qiao, Ab-
 791 hishek Nayyar, Julian Walker, Paramjit Sandhu, Mihaela Rosca, Danny Swisher, Mikhail Dek-
 792 tiarev, Josh Dillon, George-Cristian Muraru, Manuel Tragut, Artiom Myaskovsky, David Reid,
 793 Marko Velic, Owen Xiao, Jasmine George, Mark Brand, Jing Li, Wenhao Yu, Shane Gu, Xiang
 794 Deng, François-Xavier Aubet, Soheil Hassas Yeganeh, Fred Alcober, Celine Smith, Trevor Cohn,
 795 Kay McKinney, Michael Tschannen, Ramesh Sampath, Gowoon Cheon, Liangchen Luo, Luyang
 796 Liu, Jordi Orbay, Hui Peng, Gabriela Botea, Xiaofan Zhang, Charles Yoon, Cesar Magalhaes,
 797 Paweł Stradomski, Ian Mackinnon, Steven Hemingray, Kumaran Venkatesan, Rhys May, Jaeyoun
 798 Kim, Alex Druinsky, Jingchen Ye, Zheng Xu, Terry Huang, Jad Al Abdallah, Adil Dostmohamed,
 799 Rachana Fellingner, Tsendsuren Munkhdalai, Akanksha Maurya, Peter Garst, Yin Zhang, Maxim
 800 Krikun, Simon Bucher, Aditya Srikanth Veerubhotla, Yaxin Liu, Sheng Li, Nishesh Gupta, Jakub
 801 Adamek, Hanwen Chen, Bernett Orlando, Aleksandr Zaks, Joost van Amersfoort, Josh Camp, Hui
 802 Wan, HyunJeong Choe, Zhichun Wu, Kate Olszewska, Weiren Yu, Archita Vadali, Martin Scholz,
 803 Daniel De Freitas, Jason Lin, Amy Hua, Xin Liu, Frank Ding, Yichao Zhou, Boone Severson, Ka-
 804 terina Tshilas, Samuel Yang, Tammo Spalink, Varun Yerram, Helena Pankov, Rory Blevins, Ben
 805 Vargas, Sarthak Jauhari, Matt Miecnikowski, Ming Zhang, Sandeep Kumar, Clement Farabet,
 806 Charline Le Lan, Sebastian Flennerhag, Yonatan Bitton, Ada Ma, Arthur Bražinskas, Eli Collins,
 807 Niharika Ahuja, Sneha Kudugunta, Anna Bortsova, Minh Giang, Wanzheng Zhu, Ed Chi, Scott
 808 Lundberg, Alexey Stern, Subha Puttagunta, Jing Xiong, Xiao Wu, Yash Pande, Amit Jhinal,
 809 Daniel Murphy, Jon Clark, Marc Brockschmidt, Maxine Deines, Kevin R. McKee, Dan Bahir,
 Jiajun Shen, Minh Truong, Daniel McDuff, Andrea Gesmundo, Edouard Rosseel, Bowen Liang,
 Ken Caluwaerts, Jessica Hamrick, Joseph Kready, Mary Cassin, Rishikesh Ingale, Li Lao, Scott
 Pollom, Yifan Ding, Wei He, Lizzeth Bellot, Joana Iljazi, Ramya Sree Boppana, Shan Han, Tara
 Thompson, Amr Khalifa, Anna Bulanova, Blagoj Mitrevski, Bo Pang, Emma Cooney, Tian Shi,
 Rey Coaguila, Tamar Yakar, Marc’aurelio Ranzato, Nikola Momchev, Chris Rawles, Zachary

810 Charles, Young Maeng, Yuan Zhang, Rishabh Bansal, Xiaokai Zhao, Brian Albert, Yuan Yuan,
 811 Sudheendra Vijayanarasimhan, Roy Hirsch, Vinay Ramasesh, Kiran Vodrahalli, Xingyu Wang,
 812 Arushi Gupta, DJ Strouse, Jianmo Ni, Roma Patel, Gabe Taubman, Zhouyuan Huo, Dero Gharib-
 813 ian, Marianne Monteiro, Hoi Lam, Shobha Vasudevan, Aditi Chaudhary, Isabela Albuquerque,
 814 Kilol Gupta, Sebastian Riedel, Chaitra Hegde, Avraham Ruderman, András György, Marcus
 815 Wainwright, Ashwin Chaugule, Burcu Karagol Ayan, Tomer Levinboim, Sam Shleifer, Yogesh
 816 Kalley, Vahab Mirrokni, Abhishek Rao, Prabakar Radhakrishnan, Jay Hartford, Jialin Wu, Zhen-
 817 hai Zhu, Francesco Bertolini, Hao Xiong, Nicolas Serrano, Hamish Tomlinson, Myle Ott, Yifan
 818 Chang, Mark Graham, Jian Li, Marco Liang, Xiangzhu Long, Sebastian Borgeaud, Yanif Ahmad,
 819 Alex Grills, Diana Mincu, Martin Izzard, Yuan Liu, Jinyu Xie, Louis O’Bryan, Sameera Ponda,
 820 Simon Tong, Michelle Liu, Dan Malkin, Khalid Salama, Yuankai Chen, Rohan Anil, Anand Rao,
 821 Rigel Swavely, Misha Bilenko, Nina Anderson, Tat Tan, Jing Xie, Xing Wu, Lijun Yu, Oriol
 822 Vinyals, Andrey Ryabtsev, Rumén Dangovski, Kate Baumli, Daniel Keysers, Christian Wright,
 823 Zoe Ashwood, Betty Chan, Artem Shtefan, Yaohui Guo, Ankur Bapna, Radu Soricut, Steven
 824 Pecht, Sabela Ramos, Rui Wang, Jiahao Cai, Trieu Trinh, Paul Barham, Linda Friso, Eli Stickgold,
 825 Xiangzhuo Ding, Siamak Shakeri, Diego Ardila, Eleftheria Briakou, Phil Culliton, Adam Raveret,
 826 Jingyu Cui, David Saxton, Subhrajit Roy, Javad Azizi, Pengcheng Yin, Lucia Loher, Andrew Bun-
 827 ner, Min Choi, Faruk Ahmed, Eric Li, Yin Li, Shengyang Dai, Michael Elabd, Sriram Ganapathy,
 828 Shivani Agrawal, Yiqing Hua, Paige Kunkle, Sujeevan Rajayogam, Arun Ahuja, Arthur Conmy,
 829 Alex Vasiloff, Parker Beak, Christopher Yew, Jayaram Mudigonda, Bartek Wydrowski, Jon Blan-
 830 ton, Zhengdong Wang, Yann Dauphin, Zhuo Xu, Martin Polacek, Xi Chen, Hexiang Hu, Pauline
 831 Sho, Markus Kunesch, Mehdi Hafezi Manshadi, Eliza Rutherford, Bo Li, Sissie Hsiao, Iain Barr,
 832 Alex Tudor, Matija Kecman, Arsha Nagrani, Vladimir Pchelin, Martin Sundermeyer, Aishwarya P
 833 S, Abhijit Karmarkar, Yi Gao, Grishma Chole, Olivier Bachem, Isabel Gao, Arturo BC, Matt
 834 Dibb, Mauro Verzetti, Felix Hernandez-Campos, Yana Lunts, Matthew Johnson, Julia Di Trapani,
 835 Raphael Koster, Idan Brusilovsky, Binbin Xiong, Megha Mohabey, Han Ke, Joe Zou, Tea Sabolić,
 836 Víctor Campos, John Palowitch, Alex Morris, Linhai Qiu, Pranavaraj Ponnuramu, Fangtao Li,
 837 Vivek Sharma, Kiranbir Sodhia, Kaan Tekelioglu, Aleksandr Chuklin, Madhavi Yenugula, Erika
 838 Gemzer, Theofilos Strinopoulos, Sam El-Husseini, Huiyu Wang, Yan Zhong, Edouard Leurent,
 839 Paul Natsev, Weijun Wang, Dre Mahaarachchi, Tao Zhu, Songyou Peng, Sami Alabed, Cheng-
 840 Chun Lee, Anthony Brohan, Arthur Szlam, GS Oh, Anton Kovsharov, Jenny Lee, Renee Wong,
 841 Megan Barnes, Gregory Thornton, Felix Gimeno, Omer Levy, Martin Sevenich, Melvin Johnson,
 842 Jonathan Mallinson, Robert Dadashi, Ziyue Wang, Qingchun Ren, Preethi Lahoti, Arka Dhar,
 843 Josh Feldman, Dan Zheng, Thatcher Ulrich, Liviu Panait, Michiel Blokzijl, Cip Baetu, Josip
 844 Matak, Jitendra Harlalka, Maulik Shah, Tal Marian, Daniel von Dincklage, Cosmo Du, Ruy Ley-
 845 Wild, Bethanie Brownfield, Max Schumacher, Yury Stuken, Shadi Noghabi, Sonal Gupta, Xiaoqi
 846 Ren, Eric Malmi, Felix Weissenberger, Blanca Huergo, Maria Bauza, Thomas Lampe, Arthur
 847 Douillard, Mojtaba Seyedhosseini, Roy Frostig, Zoubin Ghahramani, Kelvin Nguyen, Kashyap
 848 Krishnakumar, Chengxi Ye, Rahul Gupta, Alireza Nazari, Robert Geirhos, Pete Shaw, Ahmed
 849 Eleryan, Dima Damen, Jennimaria Palomaki, Ted Xiao, Qiyin Wu, Quan Yuan, Phoenix Mead-
 850 owlark, Matthew Bilotti, Raymond Lin, Mukund Sridhar, Yannick Schroecker, Da-Woon Chung,
 851 Jincheng Luo, Trevor Strohman, Tianlin Liu, Anne Zheng, Jesse Emond, Wei Wang, Andrew
 852 Lampinen, Toshiyuki Fukuzawa, Folaayo Campbell-Ajala, Monica Roy, James Lee-Thorp, Lily
 853 Wang, Iftekhar Naim, Tony, Nguy ên, Guy Bensusky, Aditya Gupta, Dominika Rogozińska, Justin
 854 Fu, Thanumalayan Sankaranarayanan Pillai, Petar Veličković, Shahar Drath, Philipp Neubeck,
 855 Vaibhav Tulsyan, Arseniy Klimovskiy, Don Metzler, Sage Stevens, Angel Yeh, Junwei Yuan,
 856 Tianhe Yu, Kelvin Zhang, Alec Go, Vincent Tsang, Ying Xu, Andy Wan, Isaac Galatzer-Levy,
 857 Sam Sobell, Abodunrinwa Toki, Elizabeth Salesky, Wenlei Zhou, Diego Antognini, Sholto Dou-
 858 glas, Shimu Wu, Adam Lelkes, Frank Kim, Paul Cavallaro, Ana Salazar, Yuchi Liu, James Besley,
 859 Tiziana Refice, Yiling Jia, Zhang Li, Michal Sokolik, Arvind Kannan, Jon Simon, Jo Chick, Avia
 860 Aharon, Meet Gandhi, Mayank Daswani, Keyvan Amiri, Vighnesh Birodkar, Abe Ittycheriah,
 861 Peter Grabowski, Oscar Chang, Charles Sutton, Zhixin, Lai, Umesh Telang, Susie Sargsyan, Tao
 862 Jiang, Raphael Hoffmann, Nicole Brichtova, Matteo Hessel, Jonathan Halcrow, Sammy Jerome,
 863 Geoff Brown, Alex Tomala, Elena Buchatskaya, Dian Yu, Sachit Menon, Pol Moreno, Yuguo
 Liao, Vicky Zayats, Luming Tang, SQ Mah, Ashish Shenoy, Alex Siegman, Majid Hadian, Okwan
 Kwon, Tao Tu, Nima Khajehnouri, Ryan Foley, Parisa Haghani, Zhongru Wu, Vaishakh Keshava,
 Khyatti Gupta, Tony Bruguier, Rui Yao, Danny Karmon, Luisa Zintgraf, Zhicheng Wang, En-
 rique Piqueras, Junehyuk Jung, Jenny Brennan, Diego Machado, Marissa Giustina, MH Tessler,
 Kamyu Lee, Qiao Zhang, Joss Moore, Kaspar Daugaard, Alexander Frömmgen, Jennifer Beat-

864 tie, Fred Zhang, Daniel Kasenberg, Ty Geri, Danfeng Qin, Gaurav Singh Tomar, Tom Ouyang,
865 Tianli Yu, Luowei Zhou, Rajiv Mathews, Andy Davis, Yaoyiran Li, Jai Gupta, Damion Yates,
866 Linda Deng, Elizabeth Kemp, Ga-Young Joung, Sergei Vassilvitskii, Mandy Guo, Pallavi LV,
867 Dave Dopson, Sami Lachgar, Lara McConnaughey, Himadri Choudhury, Dragos Dena, Aaron
868 Cohen, Joshua Ainslie, Sergey Levi, Parthasarathy Gopavarapu, Polina Zablotskaia, Hugo Val-
869 let, Sanaz Bahargam, Xiaodan Tang, Nenad Tomasev, Ethan Dyer, Daniel Balle, Hongrae Lee,
870 William Bono, Jorge Gonzalez Mendez, Vadim Zubov, Shentao Yang, Ivor Rendulic, Yanyan
871 Zheng, Andrew Hogue, Golan Pundak, Ralph Leith, Avishkar Bhoopchand, Michael Han, Mis-
872 lav Žanić, Tom Schaul, Manolis Delakis, Tejas Iyer, Guanyu Wang, Harman Singh, Abdelrah-
873 man Abdelhamed, Tara Thomas, Siddhartha Brahma, Hilal Dib, Naveen Kumar, Wenxuan Zhou,
874 Liang Bai, Pushkar Mishra, Jiao Sun, Valentin Anklin, Roykrong Sukkerd, Lauren Agubuzu, An-
875 ton Briukhov, Anmol Gulati, Maximilian Sieb, Fabio Pardo, Sara Nasso, Junquan Chen, Kexin
876 Zhu, Tiberiu Sosea, Alex Goldin, Keith Rush, Spurthi Amba Hombaiah, Andreas Noever, Al-
877 lan Zhou, Sam Haves, Mary Phuong, Jake Ades, Yi ting Chen, Lin Yang, Joseph Pagadora,
878 Stan Bileschi, Victor Cotruta, Rachel Saputro, Arijit Pramanik, Sean Ammirati, Dan Garrette,
879 Kevin Villela, Tim Blyth, Canfer Akbulut, Neha Jha, Alban Rrustemi, Arissa Wongpanich, Chi-
880 rag Nagpal, Yonghui Wu, Morgane Rivière, Sergey Kishchenko, Pranesh Srinivasan, Alice Chen,
881 Animesh Sinha, Trang Pham, Bill Jia, Tom Hennigan, Anton Bakalov, Nithya Attaluri, Drew
882 Garmon, Daniel Rodriguez, Dawid Wegner, Wenhao Jia, Evan Senter, Noah Fiedel, Denis Pe-
883 tek, Yuchuan Liu, Cassidy Hardin, Harshal Tushar Lehri, Joao Carreira, Sara Smoot, Marcel
884 Prasetya, Nami Akazawa, Anca Stefanoiu, Chia-Hua Ho, Anelia Angelova, Kate Lin, Min Kim,
885 Charles Chen, Marcin Sieniek, Alice Li, Tongfei Guo, Sorin Baltateanu, Pouya Tafti, Michael
886 Wunder, Nadav Olmert, Divyansh Shukla, Jingwei Shen, Neel Kovelamudi, Balaji Venkatraman,
887 Seth Neel, Romal Thoppilan, Jerome Connor, Frederik Benzing, Axel Stjerngren, Golnaz Ghi-
888 asi, Alex Polozov, Joshua Howland, Theophane Weber, Justin Chiu, Ganesh Poomal Girirajan,
889 Andreas Terzis, Pidong Wang, Fangda Li, Yoav Ben Shalom, Dinesh Tewari, Matthew Denton,
890 Roe Aharoni, Norbert Kalb, Heri Zhao, Junlin Zhang, Angelos Filos, Matthew Rahtz, Lalit Jain,
891 Connie Fan, Vitor Rodrigues, Ruth Wang, Richard Shin, Jacob Austin, Roman Ring, Mariella
892 Sanchez-Vargas, Mehadi Hassen, Ido Kessler, Uri Alon, Gufeng Zhang, Wenhu Chen, Yenai Ma,
893 Xiance Si, Le Hou, Azalia Mirhoseini, Marc Wilson, Geoff Bacon, Becca Roelofs, Lei Shu,
894 Gautam Vasudevan, Jonas Adler, Artur Dwornik, Tayfun Terzi, Matt Lawlor, Harry Askham,
895 Mike Bernico, Xuanyi Dong, Chris Hidey, Kevin Kilgour, Gaël Liu, Surya Bhupatiraju, Luke
896 Leonhard, Siqi Zuo, Partha Talukdar, Qing Wei, Aliaksei Severyn, Vít Listfk, Jong Lee, Aditya
897 Tripathi, SK Park, Yossi Matias, Hao Liu, Alex Ruiz, Rajesh Jayaram, Jackson Tolins, Pierre
898 Marcenac, Yiming Wang, Bryan Seybold, Henry Prior, Deepak Sharma, Jack Weber, Mikhail
899 Sirotenko, Yunhsuan Sung, Dayou Du, Ellie Pavlick, Stefan Zinke, Markus Freitag, Max Dylla,
900 Montse Gonzalez Arenas, Natan Potikha, Omer Goldman, Connie Tao, Rachita Chhaparia, Maria
901 Voitovich, Pawan Dogra, Andrija Ražnatović, Zak Tsai, Chong You, Oleaser Johnson, George
902 Tucker, Chenjie Gu, Jae Yoo, Maryam Majzoubi, Valentin Gabeur, Bahram Raad, Rocky Rhodes,
903 Kashyap Kolipaka, Heidi Howard, Geta Sampemane, Benny Li, Chulayuth Asawaroengchai, Duy
904 Nguyen, Chiyuan Zhang, Timothee Cour, Xinxin Yu, Zhao Fu, Joe Jiang, Po-Sen Huang, Gabriela
905 Surita, Iñaki Iturrate, Yael Karov, Michael Collins, Martin Baeuml, Fabian Fuchs, Shilpa Shetty,
906 Swaroop Ramaswamy, Sayna Ebrahimi, Qiuchen Guo, Jeremy Shar, Gabe Barth-Maron, Sravanti
907 Addepalli, Bryan Richter, Chin-Yi Cheng, Eugénie Rives, Fei Zheng, Johannes Griesser, Nishanth
908 Dikkala, Yoel Zeldes, Ilkin Safarli, Dipanjan Das, Himanshu Srivastava, Sadh MNM Khan, Xin
909 Li, Aditya Pandey, Larisa Markeeva, Dan Belov, Qiqi Yan, Mikołaj Rybiński, Tao Chen, Megha
910 Nawhal, Michael Quinn, Vineetha Govindaraj, Sarah York, Reed Roberts, Roopal Garg, Namrata
911 Godbole, Jake Abernethy, Anil Das, Lam Nguyen Thiet, Jonathan Tompson, John Nham, Neera
912 Vats, Ben Caine, Wesley Helmholtz, Francesco Pongetti, Yeongil Ko, James An, Clara Huiyi Hu,
913 Yu-Cheng Ling, Julia Pawar, Robert Leland, Keisuke Kinoshita, Waleed Khawaja, Marco Selvi,
914 Eugene Ie, Danila Sinopalnikov, Lev Proleev, Nilesh Tripuraneni, Michele Bevilacqua, Seungji
915 Lee, Clayton Sanford, Dan Suh, Dustin Tran, Jeff Dean, Simon Baumgartner, Jens Heitkaemper,
916 Sagar Gubbi, Kristina Toutanova, Yichong Xu, Chandu Thekkath, Keran Rong, Palak Jain, Annie
917 Xie, Yan Virin, Yang Li, Lubo Litchev, Richard Powell, Tarun Bharti, Adam Kraft, Nan Hua,
Marissa Ikonmidis, Ayal Hitron, Sanjiv Kumar, Loic Matthey, Sophie Bridgers, Lauren Lax,
Ishaan Malhi, Ondrej Skopek, Ashish Gupta, Jiawei Cao, Michelle Rasquinha, Siim Pöder, Woj-
ciech Stokowiec, Nicholas Roth, Guowang Li, Michaël Sander, Joshua Kessinger, Vihan Jain, Ed-
ward Loper, Wonpyo Park, Michal Yarom, Liqun Cheng, Guru Guruganesh, Kanishka Rao, Yan
Li, Catarina Barros, Mikhail Sushkov, Chun-Sung Ferng, Rohin Shah, Ophir Aharoni, Ravin Ku-

918 mar, Tim McConnell, Peiran Li, Chen Wang, Fernando Pereira, Craig Swanson, Fayaz Jamil, Yan
919 Xiong, Anitha Vijayakumar, Prakash Shroff, Kedar Soparkar, Jindong Gu, Livio Baldini Soares,
920 Eric Wang, Kushal Majmundar, Aurora Wei, Kai Bailey, Nora Kassner, Chizu Kawamoto, Goran
921 Žužić, Victor Gomes, Abhirut Gupta, Michael Guzman, Ishita Dasgupta, Xinyi Bai, Zhufeng Pan,
922 Francesco Piccinno, Hadas Natalie Vogel, Octavio Ponce, Adrian Hutter, Paul Chang, Pan-Pan
923 Jiang, Ionel Gog, Vlad Ionescu, James Manyika, Fabian Pedregosa, Harry Ragan, Zach Behrman,
924 Ryan Mullins, Coline Devin, Aroonlok Pyne, Swapnil Gawde, Martin Chadwick, Yiming Gu,
925 Sasan Tavakkol, Andy Twigg, Naman Goyal, Ndidi Elue, Anna Goldie, Srinivasan Venkatachary,
926 Hongliang Fei, Ziqiang Feng, Marvin Ritter, Isabel Leal, Sudeep Dasari, Pei Sun, Alif Raditya
927 Rochman, Brendan O'Donoghue, Yuchen Liu, Jim Sproch, Kai Chen, Natalie Clay, Slav Petrov,
928 Sailesh Sidhwani, Ioana Mihailescu, Alex Panagopoulos, AJ Piergiovanni, Yunfei Bai, George
929 Powell, Deep Karkhanis, Trevor Yacovone, Petr Mitrichev, Joe Kovac, Dave Uthus, Amir Yaz-
930 danbakhsh, David Amos, Steven Zheng, Bing Zhang, Jin Miao, Bhuvana Ramabhadran, Soroush
931 Radpour, Shantanu Thakoor, Josh Newlan, Oran Lang, Orion Jankowski, Shikhar Bharadwaj,
932 Jean-Michel Sarr, Shereen Ashraf, Sneha Mondal, Jun Yan, Ankit Singh Rawat, Sarmishta Velury,
933 Greg Kochanski, Tom Eccles, Franz Och, Abhanshu Sharma, Ethan Mahintorabi, Alex Gurney,
934 Carrie Muir, Vered Cohen, Saksham Thakur, Adam Bloniarz, Asier Mujika, Alexander Pritzel,
935 Paul Caron, Altaf Rahman, Fiona Lang, Yasumasa Onoe, Petar Sirkovic, Jay Hoover, Ying Jian,
936 Pablo Duque, Arun Narayanan, David Soergel, Alex Haig, Loren Maggiore, Shyamal Buch,
937 Josef Dean, Ilya Figotin, Igor Karpov, Shaleen Gupta, Denny Zhou, Muhuan Huang, Ashwin
938 Vaswani, Christopher Semturs, Kaushik Shivakumar, Yu Watanabe, Vinodh Kumar Rajendran,
939 Eva Lu, Yanhan Hou, Wenting Ye, Shikhar Vashishth, Nana Nti, Vytenis Sakenas, Darren Ni,
940 Doug DeCarlo, Michael Bendersky, Sumit Bagri, Nacho Cano, Elijah Peake, Simon Tokumine,
941 Varun Godbole, Carlos Guía, Tanya Lando, Vittorio Selo, Seher Ellis, Danny Tarlow, Daniel
942 Gillick, Alessandro Epasto, Siddhartha Reddy Jonnalagadda, Meng Wei, Meiyang Xie, Ankur
943 Taly, Michela Paganini, Mukund Sundararajan, Daniel Toyama, Ting Yu, Dessie Petrova, Aneesh
944 Pappu, Rohan Agrawal, Senaka Buthpitiya, Justin Frye, Thomas Buschmann, Remi Crocker,
945 Marco Tagliasacchi, Mengchao Wang, Da Huang, Sagi Perel, Brian Wieder, Hideto Kazawa,
946 Weiyue Wang, Jeremy Cole, Himanshu Gupta, Ben Golan, Seojin Bang, Nitish Kulkarni, Ken
947 Franko, Casper Liu, Doug Reid, Sid Dalmia, Jay Whang, Kevin Cen, Prasha Sundaram, Johan
948 Ferret, Berivan Isik, Lucian Ionita, Guan Sun, Anna Shekhawat, Muqthar Mohammad, Philip
949 Pham, Ronny Huang, Karthik Raman, Xingyi Zhou, Ross Mcilroy, Austin Myers, Sheng Peng,
950 Jacob Scott, Paul Covington, Sofia Erell, Pratik Joshi, João Gabriel Oliveira, Natasha Noy, Tajwar
951 Nasir, Jake Walker, Vera Axelrod, Tim Dozat, Pu Han, Chun-Te Chu, Eugene Weinstein, Anand
952 Shukla, Shreyas Chandrakaladharan, Petra Poklukar, Bonnie Li, Ye Jin, Prem Eruvbetine, Steven
953 Hansen, Avigail Dabush, Alon Jacovi, Samrat Phatale, Chen Zhu, Steven Baker, Mo Shomrat,
954 Yang Xiao, Jean Pouget-Abadie, Mingyang Zhang, Fanny Wei, Yang Song, Helen King, Yiling
955 Huang, Yun Zhu, Ruoxi Sun, Juliana Vicente Franco, Chu-Cheng Lin, Sho Arora, Hui, Li, Vi-
956 vian Xia, Luke Vilnis, Mariano Schain, Kaiz Alarakyia, Laurel Prince, Aaron Phillips, Caleb
957 Habtegebriel, Luyao Xu, Huan Gui, Santiago Ontanon, Lora Aroyo, Karan Gill, Peggy Lu,
958 Yash Katariya, Dhruv Madeka, Shankar Krishnan, Shubha Srinivas Raghvendra, James Freed-
959 man, Yi Tay, Gaurav Menghani, Peter Choy, Nishita Shetty, Dan Abolafia, Doron Kukliansky,
960 Edward Chou, Jared Lichtarge, Ken Burke, Ben Coleman, Dee Guo, Larry Jin, Indro Bhat-
961 tacharya, Victoria Langston, Yiming Li, Suyog Kotecha, Alex Yakubovich, Xinyun Chen, Pe-
962 tre Petrov, Tolly Powell, Yanzhang He, Corbin Quick, Kanav Garg, Dawsen Hwang, Yang Lu,
963 Srinadh Bhojanapalli, Kristian Kjemis, Ramin Mehran, Aaron Archer, Hado van Hasselt, Ash-
964 win Balakrishna, JK Kearns, Meiqi Guo, Jason Riesa, Mikita Sazanovich, Xu Gao, Chris Sauer,
965 Chengrun Yang, XiangHai Sheng, Thomas Jimma, Wouter Van Gansbeke, Vitaly Nikolaev, Wei
966 Wei, Katie Millican, Ruizhe Zhao, Justin Snyder, Levent Bolelli, Maura O'Brien, Shawn Xu, Fei
967 Xia, Wentao Yuan, Arvind Neelakantan, David Barker, Sachin Yadav, Hannah Kirkwood, Fa-
968 rooq Ahmad, Joel Wee, Jordan Grimstad, Boyu Wang, Matthew Wiethoff, Shane Settle, Miaosen
969 Wang, Charles Blundell, Jingjing Chen, Chris Duvarney, Grace Hu, Olaf Ronneberger, Alex Lee,
970 Yuanzhen Li, Abhishek Chakladar, Alena Butryna, Georgios Evangelopoulos, Guillaume Des-
971 jardins, Jonni Kanerva, Henry Wang, Averi Nowak, Nick Li, Alyssa Loo, Art Khurshudov, Lau-
972 rent El Shafey, Nagabhushan Baddi, Karel Lenc, Yasaman Razeghi, Tom Lieber, Amer Sinha,
973 Xiao Ma, Yao Su, James Huang, Asahi Ushio, Hanna Klimczak-Plucińska, Kareem Mohamed,
974 JD Chen, Simon Osindero, Stav Ginzburg, Lampros Lamprou, Vasilisa Bashlovkina, Duc-Hieu
975 Tran, Ali Khodaei, Ankit Anand, Yixian Di, Ramy Eskander, Manish Reddy Vuyyuru, Jasmine
976 Liu, Aishwarya Kamath, Roman Goldenberg, Mathias Bellaïche, Juliette Pluto, Bill Rosgen, Has-

972 san Mansoor, William Wong, Suhas Ganesh, Eric Bailey, Scott Baird, Dan Deutsch, Jinoo Baek,
 973 Xuhui Jia, Chansoo Lee, Abe Friesen, Nathaniel Braun, Kate Lee, Amayika Panda, Steven M.
 974 Hernandez, Duncan Williams, Jianqiao Liu, Ethan Liang, Arnaud Autef, Emily Pitler, Deepali
 975 Jain, Phoebe Kirk, Oskar Bunyan, Jaume Sanchez Elias, Tongxin Yin, Machel Reid, Aedan
 976 Pope, Nikita Putikhin, Bidisha Samanta, Sergio Guadarrama, Dahun Kim, Simon Rowe, Mar-
 977 cella Valentine, Geng Yan, Alex Salcianu, David Silver, Gan Song, Richa Singh, Shuai Ye, Han-
 978 nah DeBalsi, Majd Al Merey, Eran Ofek, Albert Webson, Shibl Mourad, Ashwin Kakarla, Silvio
 979 Lattanzi, Nick Roy, Evgeny Sluzhaev, Christina Butterfield, Alessio Tonioni, Nathan Waters,
 980 Sudhindra Kopalle, Jason Chase, James Cohan, Girish Ramchandra Rao, Robert Berry, Michael
 981 Voznesensky, Shuguang Hu, Kristen Chiafullo, Sharat Chikkerur, George Scrivener, Ivy Zheng,
 982 Jeremy Wiesner, Wolfgang Macherey, Timothy Lillicrap, Fei Liu, Brian Walker, David Welling,
 983 Elinor Davies, Yangsibo Huang, Lijie Ren, Nir Shabat, Alessandro Agostini, Mariko Iinuma,
 984 Dustin Zelle, Rohit Sathyanarayana, Andrea D’olimpio, Morgan Redshaw, Matt Ginsberg, Ash-
 985 win Murthy, Mark Geller, Tatiana Matejovicova, Ayan Chakrabarti, Ryan Julian, Christine Chan,
 986 Qiong Hu, Daniel Jarrett, Manu Agarwal, Jeshwanth Challagundla, Tao Li, Sandeep Tata, Wen
 987 Ding, Maya Meng, Zhuyun Dai, Giulia Vezzani, Shefali Garg, Jannis Bulian, Mary Jasarevic,
 988 Honglong Cai, Harish Rajamani, Adam Santoro, Florian Hartmann, Chen Liang, Bartek Perz,
 989 Apoorv Jindal, Fan Bu, Sungyong Seo, Ryan Poplin, Adrian Goedeckemeyer, Badih Ghazi,
 990 Nikhil Khadke, Leon Liu, Kevin Mather, Mingda Zhang, Ali Shah, Alex Chen, Jinliang Wei, Ke-
 991 shav Shivam, Yuan Cao, Donghyun Cho, Angelo Scorza Scarpati, Michael Moffitt, Clara Barbu,
 992 Ivan Jurin, Ming-Wei Chang, Hongbin Liu, Hao Zheng, Shachi Dave, Christine Kaeser-Chen,
 993 Xiaobin Yu, Alvin Abdagic, Lucas Gonzalez, Yanping Huang, Peilin Zhong, Cordelia Schmid,
 994 Bryce Pettrini, Alex Wertheim, Jifan Zhu, Hoang Nguyen, Kaiyang Ji, Yanqi Zhou, Tao Zhou,
 995 Fangxiaoyu Feng, Regev Cohen, David Rim, Shubham Milind Phal, Petko Georgiev, Ariel Brand,
 996 Yue Ma, Wei Li, Somit Gupta, Chao Wang, Pavel Dubov, Jean Tarbouriech, Kingshuk Majumder,
 997 Huijian Li, Norman Rink, Apurv Suman, Yang Guo, Yinghao Sun, Arun Nair, Xiaowei Xu, Mo-
 998 hamed Elhawaty, Rodrigo Cabrera, Guangxing Han, Julian Eisenschlos, Junwen Bai, Yuqi Li,
 999 Yamini Bansal, Thibault Sellam, Mina Khan, Hung Nguyen, Justin Mao-Jones, Nikos Parot-
 1000 sidis, Jake Marcus, Cindy Fan, Roland Zimmermann, Yony Kochinski, Laura Graesser, Feryal
 1001 Behbahani, Alvaro Caceres, Michael Riley, Patrick Kane, Sandra Lefdal, Rob Willoughby, Paul
 1002 Vicol, Lun Wang, Shujian Zhang, Ashleah Gill, Yu Liang, Gautam Prasad, Soroosh Mariooryad,
 1003 Mehran Kazemi, Zifeng Wang, Kritika Muralidharan, Paul Voigtlaender, Jeffrey Zhao, Huanjie
 1004 Zhou, Nina D’Souza, Aditi Mavalankar, Séb Arnold, Nick Young, Obaid Sarvana, Chace Lee,
 1005 Milad Nasr, Tingting Zou, Seokhwan Kim, Lukas Haas, Kaushal Patel, Neslihan Bulut, David
 1006 Parkinson, Courtney Biles, Dmitry Kalashnikov, Chi Ming To, Aviral Kumar, Jessica Austin, Alex
 1007 Greve, Lei Zhang, Megha Goel, Yeqing Li, Sergey Yaroshenko, Max Chang, Abhishek Jindal, Ge-
 1008 off Clark, Hagai Taitelbaum, Dale Johnson, Ofir Roval, Jeongwoo Ko, Anhad Mohananeey, Chris-
 1009 tian Schuler, Shenil Dodhia, Ruichao Li, Kazuki Osawa, Claire Cui, Peng Xu, Rushin Shah, Tao
 1010 Huang, Ela Gruzewska, Nathan Clement, Mudit Verma, Olcan Sercinoglu, Hai Qian, Viral Shah,
 1011 Masa Yamaguchi, Abhinit Modi, Takahiro Kosakai, Thomas Strohmman, Junhao Zeng, Beliz
 1012 Gunel, Jun Qian, Austin Tarango, Krzysztof Jastrzebski, Robert David, Jyn Shan, Parker Schuh,
 1013 Kunal Lad, Willi Gierke, Mukundan Madhavan, Xinyi Chen, Mark Kurzeja, Rebeca Santamaria-
 1014 Fernandez, Dawn Chen, Alexandra Cordell, Yuri Chervonyi, Frankie Garcia, Nithish Kannan,
 1015 Vincent Perot, Nan Ding, Shlomi Cohen-Ganor, Victor Lavrenko, Junru Wu, Georgie Evans,
 1016 Cicero Nogueira dos Santos, Madhavi Sewak, Ashley Brown, Andrew Hard, Joan Puigcerver,
 1017 Zeyu Zheng, Yizhong Liang, Evgeny Gladchenko, Reeve Ingle, Uri First, Pierre Sermanet, Char-
 1018 lotte Magister, Mihajlo Velimirović, Sashank Reddi, Susanna Ricco, Eirikur Agustsson, Hartwig
 1019 Adam, Nir Levine, David Gaddy, Dan Holtmann-Rice, Xuanhui Wang, Ashutosh Sathe, Abhi-
 1020 jit Guha Roy, Blaž Bratanič, Alen Carin, Harsh Mehta, Silvano Bonacina, Nicola De Cao, Mara
 1021 Finkelstein, Verena Rieser, Xinyi Wu, Florent Altché, Dylan Scandinaro, Li Li, Nino Vieillard,
 1022 Nikhil Sethi, Garrett Tanzer, Zhi Xing, Shibo Wang, Parul Bhatia, Gui Citovsky, Thomas An-
 1023 thony, Sharon Lin, Tianze Shi, Shoshana Jakobovits, Gena Gibson, Raj Apte, Lisa Lee, Mingqing
 1024 Chen, Arunkumar Byravan, Petros Maniatis, Kellie Webster, Andrew Dai, Pu-Chin Chen, Jiaqi
 1025 Pan, Asya Fadeeva, Zach Gleicher, Thang Luong, and Niket Kumar Bhumihar. Gemini 2.5: Push-
 ing the frontier with advanced reasoning, multimodality, long context, and next generation agentic
 capabilities, 2025. URL <https://arxiv.org/abs/2507.06261>.

DeepSeek. Deepseek-v3.1 release. <https://api-docs.deepseek.com/news/news250821>, 2025. Accessed: 2025-09-21.

- 1026 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
1027 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
1028 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao
1029 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
1030 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,
1031 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,
1032 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
1033 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai
1034 Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang,
1035 Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang,
1036 Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang,
1037 Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang,
1038 R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng
1039 Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing
1040 Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjin Zhao, Wen
1041 Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong
1042 Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu,
1043 Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xi-
1044 aoshan Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia
1045 Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng
1046 Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong
1047 Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong,
1048 Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou,
1049 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying
1050 Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda
1051 Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu,
1052 Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu
1053 Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforce-
1054 ment learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 1055 Mostafa Dehghani, Yi Tay, Alexey A. Gritsenko, Zhe Zhao, Neil Houlsby, Fernando Diaz, Donald
1056 Metzler, and Oriol Vinyals. The benchmark lottery, 2021. URL <https://openreview.net/forum?id=5Str211vmr->.
- 1057 Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Car-
1058 oline Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, Olli
1059 Järviemi, Matthew Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla, Qiuyu Ren, Elizabeth
1060 Pratt, Lionel Levine, Grant Barkley, Natalie Stewart, Bogdan Grechuk, Tetiana Grechuk, Shreep-
1061 ranav Varma Enugandla, and Mark Wildon. Frontiermath: A benchmark for evaluating advanced
1062 mathematical reasoning in ai, 2025. URL <https://arxiv.org/abs/2411.04872>.
- 1063 Harvard–MIT Mathematics Tournament (HMMT). Hmmt february 2025 archive. Online Archive,
1064 February 2025. URL <https://www.hmmt.org/www/archive/282>. Accessed: 2025-09-
1065 21.
- 1066 Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob
1067 Steinhardt. Aligning ai with shared human values. *Proceedings of the International Conference
1068 on Learning Representations (ICLR)*, 2021a.
- 1069 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
1070 Steinhardt. Measuring massive multitask language understanding. *Proceedings of the Interna-
1071 tional Conference on Learning Representations (ICLR)*, 2021b.
- 1072 International Olympiad in Informatics. Ioi 2025, 2025. URL [https://ioi2025.obi.org.
1073 bo/tasks.html](https://ioi2025.obi.org.bo/tasks.html). Accessed: 2025-09-21.
- 1074 Alon Jacovi, Andrew Wang, Chris Alberti, Connie Tao, Jon Lipovetz, Kate Olszewska, Lukas Haas,
1075 Michelle Liu, Nate Keating, Adam Bloniarz, Carl Saroufim, Corey Fry, Dror Marcus, Doron
1076 Kukliansky, Gaurav Singh Tomar, James Swirhun, Jinwei Xing, Lily Wang, Madhu Gurusurthy,
1077 Michael Aaron, Moran Ambar, Rachana Fellingner, Rui Wang, Zizhao Zhang, Sasha Goldshtein,
1078 and Dipanjan Das. The facts grounding leaderboard: Benchmarking llms’ ability to ground re-
1079 sponses to long-form input, 2025. URL <https://arxiv.org/abs/2501.03200>.

- 1080 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Ar-
1081 mando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamina-
1082 tion free evaluation of large language models for code. In *The Thirteenth International Confer-
1083 ence on Learning Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=chfJJYC3iL)
1084 [chfJJYC3iL](https://openreview.net/forum?id=chfJJYC3iL).
- 1085 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
1086 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth
1087 International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=VTF8yNQm66)
1088 [net/forum?id=VTF8yNQm66](https://openreview.net/forum?id=VTF8yNQm66).
- 1089 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
1090 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebas-
1091 tian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts,
1092 and Adina Williams. Dynabench: Rethinking benchmarking in NLP. In Kristina Toutanova,
1093 Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cot-
1094 terell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of
1095 the North American Chapter of the Association for Computational Linguistics: Human Lan-
1096 guage Technologies*, pp. 4110–4124, Online, June 2021. Association for Computational Linguis-
1097 tics. doi: 10.18653/v1/2021.naacl-main.324. URL [https://aclanthology.org/2021.](https://aclanthology.org/2021.naacl-main.324/)
1098 [naacl-main.324/](https://aclanthology.org/2021.naacl-main.324/).
- 1099 Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang.
1100 Benchmarking cognitive biases in large language models as evaluators. In Lun-Wei Ku, Andre
1101 Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics:
1102 ACL 2024*, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL
1103 <https://aclanthology.org/2024.findings-acl.29/>.
- 1104 Qintong Li, Leyang Cui, Lingpeng Kong, and Wei Bi. Exploring the reliability of large lan-
1105 guage models as customized evaluators for diverse NLP tasks. In Owen Rambow, Leo Wan-
1106 ner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert (eds.),
1107 *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 10325–
1108 10344, Abu Dhabi, UAE, January 2025a. Association for Computational Linguistics. URL
1109 <https://aclanthology.org/2025.coling-main.688/>.
- 1110 Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E. Gon-
1111 zalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and
1112 benchbuilder pipeline. In *Forty-second International Conference on Machine Learning*, 2025b.
1113 URL <https://openreview.net/forum?id=KfTf9vFvSn>.
- 1114 Xiang Lisa Li, Farzaan Kaiyom, Evan Zheran Liu, Yifan Mai, Percy Liang, and Tatsunori
1115 Hashimoto. Autobench: Towards declarative benchmark construction. In *The Thirteenth In-
1116 ternational Conference on Learning Representations*, 2025c. URL [https://openreview.](https://openreview.net/forum?id=ymt4crbbXh)
1117 [net/forum?id=ymt4crbbXh](https://openreview.net/forum?id=ymt4crbbXh).
- 1118 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga,
1119 Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang
1120 Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D Manning, Christopher Re,
1121 Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda
1122 Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng,
1123 Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khat-
1124 tab, Peter Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar,
1125 Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William
1126 Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of lan-
1127 guage models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL
1128 <https://openreview.net/forum?id=iO4LZibEqw>. Featured Certification, Expert
1129 Certification, Outstanding Certification.
- 1130 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan
1131 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth
1132 International Conference on Learning Representations*, 2024. URL [https://openreview.](https://openreview.net/forum?id=v8L0pN6EOi)
1133 [net/forum?id=v8L0pN6EOi](https://openreview.net/forum?id=v8L0pN6EOi).

- 1134 Nikolai Liubimov. Everybody is (unintentionally) cheating: Fixing ai benchmarks.
 1135 Blog post, Label Studio, May 2025. URL [https://labelstud.io/blog/
 1136 everybody-is-unintentionally-cheating/](https://labelstud.io/blog/everybody-is-unintentionally-cheating/). Accessed: 2025-09-21.
 1137
- 1138 OpenAI. Introducing gpt-5. <https://openai.com/index/introducing-gpt-5/>,
 1139 2025. Accessed: 2025-09-21.
- 1140 Samuel J Paech. Eq-bench creative writing benchmark v3. [https://github.com/
 1141 EQ-bench/creative-writing-bench](https://github.com/EQ-bench/creative-writing-bench), 2025.
 1142
- 1143 Pankaj Pandey. Ai benchmarks are broken — here’s how we can fix them. Medium,
 1144 Feb 2025. URL [https://medium.com/@publicapplicationcenter/
 1145 ai-benchmarks-are-broken-heres-how-we-can-fix-them-eca9a521b649](https://medium.com/@publicapplicationcenter/ai-benchmarks-are-broken-heres-how-we-can-fix-them-eca9a521b649).
 1146 Accessed: 2025-09-21.
- 1147 Rob Petrosino. Ai benchmarks are broken: Why we need harder tests before it’s too
 1148 late. LinkedIn Pulse, Aug 2025. URL [https://www.linkedin.com/pulse/
 1149 ai-benchmarks-broken-why-we-need-harder-tests-before-its-petrosino-kzqle/](https://www.linkedin.com/pulse/ai-benchmarks-broken-why-we-need-harder-tests-before-its-petrosino-kzqle/).
 1150 Accessed: 2025-09-21.
 1151
- 1152 Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin
 1153 Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra,
 1154 Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Dmitry
 1155 Dodonov, Tung Nguyen, Jaeho Lee, Daron Anderson, Mikhail Doroshenko, Alun Cennyth
 1156 Stokes, Mobeen Mahmood, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, John-Clark Levin,
 1157 Mstyslav Kazakov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal,
 1158 Chelsea Zou, Zihan Wang, Serguei Popov, Robert Gerbicz, Geoff Galgon, Johannes Schmitt, Will
 1159 Yeadon, Yongki Lee, Scott Sauers, Alvaro Sanchez, Fabian Giska, Marc Roth, Søren Riis, Saiteja
 1160 Utpala, Noah Burns, Gashaw M. Goshu, Mohinder Maheshbhai Naiya, Chidozie Agu, Zachary
 1161 Giboney, Antrell Cheatom, Francesco Fournier-Facio, Sarah-Jane Crowson, Lennart Finke, Zerui
 1162 Cheng, Jennifer Zampese, Ryan G. Hoerr, Mark Nandor, Hyunwoo Park, Tim Gehrunger, Ji-
 1163 aqi Cai, Ben McCarty, Alexis C Garretson, Edwin Taylor, Damien Sileo, Qiuyu Ren, Usman
 1164 Qazi, Lianghui Li, Jungbae Nam, John B. Wydallis, Pavel Arkhipov, Jack Wei Lun Shi, Aras
 1165 Bacho, Chris G. Willcocks, Hangrui Cao, Sumeet Motwani, Emily de Oliveira Santos, Johannes
 1166 Veith, Edward Vendrow, Doru Cojoc, Kengo Zenitani, Joshua Robinson, Longke Tang, Yuqi Li,
 1167 Joshua Vendrow, Natanael Wildner Fraga, Vladyslav Kuchkin, Andrey Pupasov Maksimov, Pierre
 1168 Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Aleksandar Mikov, Andrew Gritsevskiy,
 1169 Julien Guillod, Gözdenur Demir, Dakotah Martinez, Ben Pageler, Kevin Zhou, Saeed Soori,
 1170 Ori Press, Henry Tang, Paolo Rissone, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric
 1171 Dieuleveut, Joseph Marvin Imperial, Ameya Prabhu, Jinzhou Yang, Nick Crispino, Arun Rao,
 1172 Dimitri Zvonkine, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate
 1173 Stambaugh, Subrata Mishra, Tad Hogg, Carlo Bosio, Brian P Coppola, Julian Salazar, Jaehyeok
 1174 Jin, Rafael Sayous, Stefan Ivanov, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran,
 1175 Andres Algaba, Kelsey Van den Houte, Lynn Van Der Sypt, Brecht Verbeken, David Noever,
 1176 Alexei Kopylov, Benjamin Myklebust, Bikun Li, Lisa Schut, Evgenii Zheltonozhskii, Qiaochu
 1177 Yuan, Derek Lim, Richard Stanley, Tong Yang, John Maar, Julian Wykowski, Martí Oller, An-
 1178 mol Sahu, Cesare Giulio Ardito, Yuzheng Hu, Ariel Ghislain Kemogne Kamdoum, Alvin Jin,
 1179 Tobias Garcia Vilchis, Yuexuan Zu, Martin Lackner, James Koppel, Gongbo Sun, Daniil S. Anto-
 1180 nenko, Steffi Chern, Bingchen Zhao, Pierrot Arsene, Joseph M Cavanagh, Daofeng Li, Jiawei
 1181 Shen, Donato Crisostomi, Wenjin Zhang, Ali Dehghan, Sergey Ivanov, David Perrella, Nur-
 1182 din Kaparov, Allen Zang, Iliia Sucholutsky, Arina Kharlamova, Daniil Orel, Vladislav Porit-
 1183 ski, Shalev Ben-David, Zachary Berger, Parker Whitfill, Michael Foster, Daniel Munro, Linh
 1184 Ho, Shankar Sivarajan, Dan Bar Hava, Aleksey Kuchkin, David Holmes, Alexandra Rodriguez-
 1185 Romero, Frank Sommerhage, Anji Zhang, Richard Moat, Keith Schneider, Zakayo Kazibwe,
 1186 Don Clarke, Dae Hyun Kim, Felipe Meneguitti Dias, Sara Fish, Veit Elser, Tobias Kreiman, Vic-
 1187 tor Efrén Guadarrama Vilchis, Immo Klose, Ujjwala Anantheswaran, Adam Zweiger, Kaivalya
 Rawal, Jeffery Li, Jeremy Nguyen, Nicolas Daans, Haline Heidinger, Maksim Radionov, Václav
 Rozhoň, Vincent Ginis, Christian Stump, Niv Cohen, Rafał Poświata, Josef Tkadlec, Alan Gold-
 farb, Chenguang Wang, Piotr Padlewski, Stanislaw Barzowski, Kyle Montgomery, Ryan Stendall,
 Jamie Tucker-Foltz, Jack Stade, T. Ryan Rogers, Tom Goertzen, Declan Grabb, Abhishek Shukla,

1188 Alan Givré, John Arnold Ambay, Archan Sen, Muhammad Fayez Aziz, Mark H Inlow, Hao He,
 1189 Ling Zhang, Younesse Kaddar, Ivar Ångquist, Yanxu Chen, Harrison K Wang, Kalyan Ramakr-
 1190 ishnan, Elliott Thornley, Antonio Terpin, Hailey Schoelkopf, Eric Zheng, Avisly Carmi, Ethan
 1191 D. L. Brown, Kelin Zhu, Max Bartolo, Richard Wheeler, Martin Stehberger, Peter Bradshaw,
 1192 JP Heimonen, Kaustubh Sridhar, Ido Akov, Jennifer Sandlin, Yury Makarychev, Joanna Tam,
 1193 Hieu Hoang, David M. Cunningham, Vladimir Goryachev, Demosthenes Patramanis, Michael
 1194 Krause, Andrew Redenti, David Aldous, Jesyin Lai, Shannon Coleman, Jiangnan Xu, Sang-
 1195 won Lee, Ilias Magoulas, Sandy Zhao, Ning Tang, Michael K. Cohen, Orr Paradise, Jan Hen-
 1196 drik Kirchner, Maksym Ovchynnikov, Jason O. Matos, Adithya Shenoy, Michael Wang, Yuzhou
 1197 Nie, Anna Szyber-Betley, Paolo Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Halasyamani,
 1198 Shreyas Verma, Prashant Joshi, Eli Meril, Ziqiao Ma, Jérémy Andréoletti, Raghav Singhal, Jacob
 1199 Platnick, Volodymyr Nevirkovets, Luke Basler, Alexander Ivanov, Seri Khoury, Nils Gustafsson,
 1200 Marco Piccardo, Hamid Mostaghimi, Qijia Chen, Virendra Singh, Tran Quoc Khánh, Paul Rosu,
 1201 Hannah Szlyk, Zachary Brown, Himanshu Narayan, Aline Menezes, Jonathan Roberts, William
 1202 Alley, Kunyang Sun, Arkil Patel, Max Lamparth, Anka Reuel, Linwei Xin, Hanmeng Xu, Jacob
 1203 Loader, Freddie Martin, Zixuan Wang, Andrea Achilleos, Thomas Preu, Tomek Korbak, Ida Bo-
 1204 sio, Fereshteh Kazemi, Ziye Chen, Biró Bálint, Eve J. Y. Lo, Jiaqi Wang, Maria Inês S. Nunes,
 1205 Jeremiah Milbauer, M Saiful Bari, Zihao Wang, Behzad Ansarinejad, Yewen Sun, Stephane Du-
 1206 rand, Hossam Elgnainy, Guillaume Douville, Daniel Tordera, George Balabanian, Hew Wolff,
 1207 Lynna Kvistad, Hsiaoyun Milliron, Ahmad Sakor, Murat Eron, Andrew Favre D. O., Shailesh
 1208 Shah, Xiaoxiang Zhou, Firuz Kamalov, Sherwin Abdoli, Tim Santens, Shaul Barkan, Allison
 1209 Tee, Robin Zhang, Alessandro Tomasiello, G. Bruno De Luca, Shi-Zhuo Looi, Vinh-Kha Le,
 1210 Noam Kolt, Jiayi Pan, Emma Rodman, Jacob Drori, Carl J Fossum, Niklas Muennighoff, Milind
 1211 Jagota, Ronak Pradeep, Honglu Fan, Jonathan Eicher, Michael Chen, Kushal Thaman, William
 1212 Merrill, Moritz Firsching, Carter Harris, Stefan Ciobăcă, Jason Gross, Rohan Pandey, Ilya Gusev,
 1213 Adam Jones, Shashank Agnihotri, Pavel Zhelnov, Mohammadreza Mofayezi, Alexander Piper-
 1214 ski, David K. Zhang, Kostiantyn Dobarskyi, Roman Leventov, Ignat Soroko, Joshua Duersch,
 1215 Vage Taamazyan, Andrew Ho, Wenjie Ma, William Held, Ruicheng Xian, Armel Randy Ze-
 1216 baze, Mohanad Mohamed, Julian Noah Leser, Michelle X Yuan, Laila Yacar, Johannes Lengler,
 1217 Katarzyna Olszewska, Claudio Di Fratta, Edson Oliveira, Joseph W. Jackson, Andy Zou, Muthu
 1218 Chidambaram, Timothy Manik, Hector Haffenden, Dashiell Stander, Ali Dasouqi, Alexander
 1219 Shen, Bitu Golshani, David Stap, Egor Kretov, Mikalai Uzhou, Alina Borisovna Zhidkovskaya,
 1220 Nick Winter, Miguel Orbegozo Rodriguez, Robert Lauff, Dustin Wehr, Colin Tang, Zaki Hos-
 1221 sain, Shaun Phillips, Fortuna Samuele, Fredrik Ekström, Angela Hammon, Oam Patel, Faraz
 1222 Farhidi, George Medley, Forough Mohammadzadeh, Madellene Peñaflo, Haile Kassahun, Alena
 1223 Friedrich, Rayner Hernandez Perez, Daniel Pyda, Taom Sakal, Omkar Dhamane, Ali Khajegili
 1224 Mirabadi, Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mohammad Maghsoudimehrabani,
 1225 Alon Amit, Dave Hulbert, Roberto Pereira, Simon Weber, Handoko, Anton Peristy, Stephen
 1226 Malina, Mustafa Mehkary, Rami Aly, Frank Reidegeld, Anna-Katharina Dick, Cary Friday,
 1227 Mukhwinder Singh, Hassan Shapourian, Wanyoung Kim, Mariana Costa, Hubeyb Gurdogan,
 1228 Harsh Kumar, Chiara Ceconello, Chao Zhuang, Haon Park, Micah Carroll, Andrew R. Tawfeek,
 1229 Stefan Steinerberger, Daattavya Aggarwal, Michael Kirchhof, Linjie Dai, Evan Kim, Johan Fer-
 1230 ret, Jainam Shah, Yuzhou Wang, Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Antonio Franca,
 1231 Diana T. Pham, Kang Yong Loh, Joshua Robinson, Abram Jackson, Paolo Giordano, Philipp
 1232 Petersen, Adrian Cosma, Jesus Colino, Colin White, Jacob Votava, Vladimir Vinnikov, Ethan
 1233 Delaney, Petr Spelda, Vit Stritecky, Syed M. Shahid, Jean-Christophe Mourrat, Lavr Vetoshkin,
 1234 Koen Sponselee, Renas Bacho, Zheng-Xin Yong, Florencia de la Rosa, Nathan Cho, Xiuyu Li,
 1235 Guillaume Malod, Orion Weller, Guglielmo Albani, Leon Lang, Julien Laurendeau, Dmitry Kaza-
 1236 kov, Fatimah Adesanya, Julien Portier, Lawrence Hollom, Victor Souza, Yuchen Anna Zhou,
 1237 Julien Degorre, Yiğit Yalın, Gbenga Daniel Obikoya, Rai, Filippo Bigi, M. C. Boscá, Oleg Shu-
 1238 mar, Kaniuar Bacho, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie,
 1239 Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu,
 1240 Stefano Cavalleri, Olle Häggström, Emil Verkama, Joshua Newbould, Hans Gundlach, Leonor
 1241 Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Ting Wang, Yosi Kratish, Wen-Ding
 Li, Sivakanth Gopi, Andrea Caciolai, Christian Schroeder de Witt, Pablo Hernández-Cámara,
 Emanuele Rodolà, Jules Robins, Dominic Williamson, Vincent Cheng, Brad Raynor, Hao Qi,
 Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Bren-
 ner, Mao Mao, Christoph Demian, Peyman Kassani, Xinyu Zhang, David Avagian, Eshawn Jes-
 sica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk

1242 Bodur, D. P. Shinde, Yan Carlos Leyva Labrador, Zahra Adoul, Mohamed Zekry, Ali Karakoc,
 1243 Tania C. B. Santos, Samir Shamseldeen, Loukmane Karim, Anna Liakhovitskaia, Nate Resman,
 1244 Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Earth Anderson, Rodrigo De Oliveira
 1245 Pena, Elizabeth Kelley, Hodjat Mariji, Rasoul Pouriamanesh, Wentao Wu, Ross Finocchio, Is-
 1246 mail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai,
 1247 Siriphan Arthornthurasuk, Isaac C. McAlister, Alejandro José Moyano, Alexey Pronin, Jing Fan,
 1248 Angel Ramirez-Trinidad, Yana Malysheva, Daphiny Pottmaier, Omid Taheri, Stanley Stepanic,
 1249 Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Ricardo Lorena,
 1250 Krishnamurthy Iyer, Arshad Anil Fasiludeen, Ronald Clark, Josh Ducey, Matheus Piza, Maja
 1251 Somrak, Eric Vergo, Juehang Qin, Benjámín Borbás, Eric Chu, Jack Lindsey, Antoine Jallon,
 1252 I. M. J. McInnis, Evan Chen, Avi Semler, Luk Gloor, Tej Shah, Marc Carauleanu, Pascal Lauer,
 1253 Tran Duc Huy, Hossein Shahrtash, Emilien Duc, Lukas Lewark, Assaf Brown, Samuel Albanie,
 1254 Brian Weber, Warren S. Vaz, Pierre Clavier, Yiyang Fan, Gabriel Poesia Reis e Silva, Long,
 1255 Lian, Marcus Abramovitch, Xi Jiang, Sandra Mendoza, Murat Islam, Juan Gonzalez, Vasilios
 1256 Mavroudis, Justin Xu, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari,
 1257 Ferenc Jeanplong, Thorben Jansen, Antonella Pinto, Archimedes Apronti, Abdallah Galal, Ng Ze-
 1258 An, Ankit Singh, Tong Jiang, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani,
 1259 Gang Zhang, Zhehang Du, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy,
 1260 Taylor D. Hartman, Tim Tarver, Stephen Mensah, Gautier Abou Loume, Wiktor Morak, Farzad
 1261 Habibi, Sarah Hoback, Will Cai, Javier Gimenez, Roselynn Grace Montecillo, Jakub Łucki, Rus-
 1262 sell Campbell, Asankhaya Sharma, Khalida Meer, Shreen Gul, Daniel Espinosa Gonzalez, Xavier
 1263 Alapont, Alex Hoover, Gunjan Chhablani, Freddie Vargus, Arunim Agarwal, Yibo Jiang, Deep-
 1264 akkumar Patil, David Outevsky, Kevin Joseph Scaria, Rajat Maheshwari, Abdelkader Dendane,
 1265 Priti Shukla, Ashley Cartwright, Sergei Bogdanov, Niels Müндler, Sören Möller, Luca Arnaboldi,
 1266 Kunvar Thaman, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Tony Fruhauff,
 1267 Glen Sherman, Mátyás Vincze, Siranut Usawasutsakorn, Dylan Ler, Anil Radhakrishnan, In-
 1268 nocent Enyekwe, Sk Md Salauddin, Jiang Muzhen, Aleksandr Maksapetyan, Vivien Rossbach,
 1269 Chris Harjadi, Mohsen Bahalooohoreh, Claire Sparrow, Jasdeep Sidhu, Sam Ali, Song Bian, John
 1270 Lai, Eric Singer, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Menshawy, Darling
 1271 Duclosel, Dario Bezzi, Yashaswini Jain, Ashley Aaron, Murat Tiryakioglu, Sheeshram Siddh,
 1272 Keith Krenek, Imad Ali Shah, Jun Jin, Scott Creighton, Denis Peskoff, Zienab EL-Wasif, Ra-
 1273 gavendran P V, Michael Richmond, Joseph McGowan, Tejal Patwardhan, Hao-Yu Sun, Ting Sun,
 1274 Nikola Zubić, Samuele Sala, Stephen Ebert, Jean Kaddour, Manuel Schottdorf, Dianzhuo Wang,
 1275 Gerol Petruzella, Alex Meiburg, Tilen Medved, Ali ElSheikh, S Ashwin Hebbar, Lorenzo Va-
 1276 quero, Xianjun Yang, Jason Poulos, Vilém Zouhar, Sergey Bogdanik, Mingfang Zhang, Jorge
 1277 Sanz-Ros, David Anugraha, Yinwei Dai, Anh N. Nhu, Xue Wang, Ali Anil Demircali, Zhibai Jia,
 1278 Yuyin Zhou, Juncheng Wu, Mike He, Nitin Chandok, Aarush Sinha, Gaoxiang Luo, Long Le,
 1279 Mickaël Noyé, Michał Perelkiewicz, Ioannis Pantidis, Tianbo Qi, Soham Sachin Purohit, Letitia
 1280 Parcalabescu, Thai-Hoa Nguyen, Genta Indra Winata, Edoardo M. Ponti, Hanchen Li, Kaustubh
 1281 Dhole, Jongee Park, Dario Abbondanza, Yuanli Wang, Anupam Nayak, Diogo M. Caetano, Anto-
 1282 nio A. W. L. Wong, Maria del Rio-Chanona, Dániel Kondor, Pieter Francois, Ed Chalstrey, Jakob
 1283 Zsombok, Dan Hoyer, Jenny Reddish, Jakob Hauer, Francisco-Javier Rodrigo-Ginés, Suchandra
 1284 Datta, Maxwell Shepherd, Thom Kamphuis, Qizheng Zhang, Hyunjun Kim, Ruiji Sun, Jianzhu
 1285 Yao, Franck Dernoncourt, Satyapriya Krishna, Sina Rismanchian, Bonan Pu, Francesco Pinto,
 1286 Yingheng Wang, Kumar Shridhar, Kalon J. Overholt, Glib Briia, Hieu Nguyen, David, Soler Bar-
 1287 tomeu, Tony CY Pang, Adam Wecker, Yifan Xiong, Fanfei Li, Lukas S. Huber, Joshua Jaeger,
 1288 Romano De Maddalena, Xing Han Lù, Yuhui Zhang, Claas Beger, Patrick Tser Jern Kon, Sean Li,
 1289 Vivek Sanker, Ming Yin, Yihao Liang, Xinlu Zhang, Ankit Agrawal, Li S. Yifei, Zechen Zhang,
 1290 Mu Cai, Yasin Sonmez, Costin Cozianu, Changhao Li, Alex Slen, Shoubin Yu, Hyun Kyu Park,
 1291 Gabriele Sarti, Marcin Briański, Alessandro Stolfo, Truong An Nguyen, Mike Zhang, Yotam
 1292 Perlit, Jose Hernandez-Orallo, Runjia Li, Amin Shabani, Felix Juefei-Xu, Shikhar Dhingra,
 1293 Orr Zohar, My Chiffon Nguyen, Alexander Pondaven, Abdurrahim Yilmaz, Xuandong Zhao,
 1294 Chuanyang Jin, Muyan Jiang, Stefan Todoran, Xinyao Han, Jules Kreuer, Brian Rabern, Anna
 1295 Plassart, Martino Maggetti, Luther Yap, Robert Geirhos, Jonathon Kean, Dingsu Wang, Sina
 Mollaei, Chenkai Sun, Yifan Yin, Shiqi Wang, Rui Li, Yaowen Chang, Anjiang Wei, Alice
 Bizeul, Xiaohan Wang, Alexandre Oliveira Arrais, Kushin Mukherjee, Jorge Chamorro-Padial,
 Jiachen Liu, Xingyu Qu, Junyi Guan, Adam Bouyamourn, Shuyu Wu, Martyna Plomecka, Junda
 Chen, Mengze Tang, Jiaqi Deng, Shreyas Subramanian, Haocheng Xi, Haoxuan Chen, Weizhi
 Zhang, Yinuo Ren, Haoqin Tu, Sejong Kim, Yushun Chen, Sara Vera Marjanović, Junwoo Ha,

- 1296 Grzegorz Luczyna, Jeff J. Ma, Zewen Shen, Dawn Song, Cedegao E. Zhang, Zhun Wang, Gaël
1297 Gendron, Yunze Xiao, Leo Smucker, Erica Weng, Kwok Hao Lee, Zhe Ye, Stefano Ermon, Ig-
1298 nacio D. Lopez-Miguel, Theo Knights, Anthony Gitter, Namkyu Park, Boyi Wei, Hongzheng
1299 Chen, Kunal Pai, Ahmed Elkhanany, Han Lin, Philipp D. Siedler, Jichao Fang, Ritwik Mishra,
1300 Károly Zsolnai-Fehér, Xilin Jiang, Shadab Khan, Jun Yuan, Rishab Kumar Jain, Xi Lin, Mike
1301 Peterson, Zhe Wang, Aditya Malusare, Maosen Tang, Isha Gupta, Ivan Fosin, Timothy Kang,
1302 Barbara Dworakowska, Kazuki Matsumoto, Guangyao Zheng, Gerben Sewuster, Jorge Pretel
1303 Villanueva, Ivan Rannev, Igor Chernyavsky, Jiale Chen, Deepayan Banik, Ben Racz, Wenchao
1304 Dong, Jianxin Wang, Laila Bashmal, Duarte V. Gonçalves, Wei Hu, Kaushik Bar, Ondrej Bo-
1305 hdal, Atharv Singh Patlan, Shehzaad Dhuliawala, Caroline Geirhos, Julien Wist, Yuval Kansal,
1306 Bingsen Chen, Kutay Tire, Atak Talay Yücel, Brandon Christof, Veerupaksh Singla, Zijian Song,
1307 Sanxing Chen, Jiaxin Ge, Kaustubh Ponkshe, Isaac Park, Tianneng Shi, Martin Q. Ma, Joshua
1308 Mak, Sherwin Lai, Antoine Moulin, Zhuo Cheng, Zhanda Zhu, Ziyi Zhang, Vaidehi Patil, Ketan
1309 Jha, Qitong Men, Jiaxuan Wu, Tianchi Zhang, Bruno Hebling Vieira, Alham Fikri Aji, Jae-Won
1310 Chung, Mohammed Mahfoud, Ha Thi Hoang, Marc Sperzel, Wei Hao, Kristof Meding, Sihan
1311 Xu, Vassilis Kostakos, Davide Manini, Yueying Liu, Christopher Toukmaji, Jay Paek, Eunmi Yu,
1312 Arif Engin Demircali, Zhiyi Sun, Ivan Dewerpe, Hongsen Qin, Roman Pflugfelder, James Bailey,
1313 Johnathan Morris, Ville Heilala, Sybille Rosset, Zishun Yu, Peter E. Chen, Woongyeong Yeo, Ee-
1314 shaan Jain, Ryan Yang, Sreekar Chigurupati, Julia Chernyavsky, Sai Prajwal Reddy, Subhashini
1315 Venugopalan, Hunar Batra, Core Francisco Park, Hieu Tran, Guilherme Maximiano, Genghan
1316 Zhang, Yizhuo Liang, Hu Shiyu, Rongwu Xu, Rui Pan, Siddharth Suresh, Ziqi Liu, Samaksh Gu-
1317 lati, Songyang Zhang, Peter Turchin, Christopher W. Bartlett, Christopher R. Scotese, Phuong M.
1318 Cao, Aakaash Nattanmai, Gordon McKellips, Anish Cheraku, Asim Suhail, Ethan Luo, Marvin
1319 Deng, Jason Luo, Ashley Zhang, Kavin Jindel, Jay Paek, Kasper Halevy, Allen Baranov, Michael
1320 Liu, Advait Avadhanam, David Zhang, Vincent Cheng, Brad Ma, Evan Fu, Liam Do, Joshua
1321 Lass, Hubert Yang, Surya Sunkari, Vishruth Bharath, Violet Ai, James Leung, Rishit Agrawal,
1322 Alan Zhou, Kevin Chen, Tejas Kalpathi, Ziqi Xu, Gavin Wang, Tyler Xiao, Erik Maung, Sam
1323 Lee, Ryan Yang, Roy Yue, Ben Zhao, Julia Yoon, Sunny Sun, Aryan Singh, Ethan Luo, Clark
1324 Peng, Tyler Osbey, Taozhi Wang, Daryl Echeazu, Hubert Yang, Timothy Wu, Spandan Patel,
1325 Vidhi Kulkarni, Vijaykaarti Sundarapandiyan, Ashley Zhang, Andrew Le, Zafir Nasim, Srikanth
1326 Yalam, Ritesh Kasamsetty, Soham Samal, Hubert Yang, David Sun, Nihar Shah, Abhijeet Saha,
1327 Alex Zhang, Leon Nguyen, Laasya Nagumalli, Kaixin Wang, Alan Zhou, Aidan Wu, Jason Luo,
1328 Anwith Telluri, Summer Yue, Alexandr Wang, and Dan Hendrycks. Humanity’s last exam, 2025.
1329 URL <https://arxiv.org/abs/2501.14249>.
- 1330 Valentina Pyatkin, Saumya Malik, Victoria Graf, Hamish Ivison, Shengyi Huang, Pradeep Dasigi,
1331 Nathan Lambert, and Hannaneh Hajishirzi. Generalizing verifiable instruction following, 2025.
1332 URL <https://arxiv.org/abs/2507.02833>.
- 1333 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
1334 Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a
1335 benchmark. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- 1336 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
1337 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Lan-
1338 guage models are multilingual chain-of-thought reasoners. In *The Eleventh International Confer-
1339 ence on Learning Representations*, 2023. URL <https://openreview.net/forum?id=fR3wGck-IXp>.
- 1340 Luísa Shimabucoro, Yongshuo Zong, Timothy Hospedales, and Henry Gouk. Evaluating the
1341 evaluators: Are current few-shot learning benchmarks fit for purpose?, 2024. URL <https://openreview.net/forum?id=kiwyQsZIGP>.
- 1342 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen,
1343 Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong,
1344 Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao,
1345 Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang
1346 Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu,
1347 Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin,

- 1350 Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao
1351 Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin
1352 Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu,
1353 Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe
1354 Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo
1355 Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi,
1356 Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng
1357 Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaying Wang,
1358 Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang,
1359 Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu,
1360 Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing
1361 Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie
1362 Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao,
1363 Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang
1364 Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang,
1365 Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng
1366 Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou,
1367 Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence,
2025a. URL <https://arxiv.org/abs/2507.20534>.
- 1368
1369 P Team, Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu,
1370 Yiming Liang, Xiaolong Jin, Zhenlin Wei, Chujie Zheng, Kaixin Deng, Shawn Gavin, Shian
1371 Jia, Sichao Jiang, Yiyan Liao, Rui Li, Qinrui Li, Sirun Li, Yizhi Li, Yunwen Li, David Ma,
1372 Yuansheng Ni, Haoran Que, Qiyao Wang, Zhoufutu Wen, Siwei Wu, Tyshawn Hsing, Ming
1373 Xu, Zhenzhu Yang, Zekun Moore Wang, Junting Zhou, Yuelin Bai, Xingyuan Bu, Chenglin
1374 Cai, Liang Chen, Yifan Chen, Chengtuo Cheng, Tianhao Cheng, Keyi Ding, Siming Huang,
1375 Yun Huang, Yaoru Li, Yizhe Li, Zhaoqun Li, Tianhao Liang, Chengdong Lin, Hongquan Lin,
1376 Yinghao Ma, Tianyang Pang, Zhongyuan Peng, Zifan Peng, Qige Qi, Shi Qiu, Xingwei Qu,
1377 Shanghaoran Quan, Yizhou Tan, Zili Wang, Chenqing Wang, Hao Wang, Yiya Wang, Yubo
1378 Wang, Jiajun Xu, Kexin Yang, Ruibin Yuan, Yuanhao Yue, Tianyang Zhan, Chun Zhang, Jinyang
1379 Zhang, Xiyue Zhang, Xingjian Zhang, Yue Zhang, Yongchi Zhao, Xiangyu Zheng, Chenghua
1380 Zhong, Yang Gao, Zhoujun Li, Dayiheng Liu, Qian Liu, Tianyu Liu, Shiwen Ni, Junran Peng,
1381 Yujia Qin, Wenbo Su, Guoyin Wang, Shi Wang, Jian Yang, Min Yang, Meng Cao, Xiang
1382 Yue, Zhaoxiang Zhang, Wangchunshu Zhou, Jiaheng Liu, Qunshu Lin, Wenhao Huang, and
1383 Ge Zhang. Supergpqa: Scaling llm evaluation across 285 graduate disciplines, 2025b. URL
<https://arxiv.org/abs/2502.14739>.
- 1384 The Terminal-Bench Team. Terminal-bench: A benchmark for ai agents in terminal environments,
1385 Apr 2025. URL <https://github.com/laude-institute/terminal-bench>.
- 1386
1387 Theo Von. Sam altman — this past weekend w/ theo von. [https://www.youtube.com/](https://www.youtube.com/watch?v=aYn8VKW6vXA)
1388 [watch?v=aYn8VKW6vXA](https://www.youtube.com/watch?v=aYn8VKW6vXA), 2025. Accessed: 2025-09-21.
- 1389
1390 Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland
1391 Haas, Pan Ji, Kittithat Krongchon, Yao Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha
1392 Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Minhui Zhu, Kilian
1393 Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, Eliu Huerta,
1394 and Hao Peng. Scicode: A research coding benchmark curated by scientists, 2024. URL <https://arxiv.org/abs/2407.13168>.
- 1395
1396 Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Courn-
1397 peau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der
1398 Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nel-
1399 son, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore,
1400 Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quin-
1401 tero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van
1402 Mulbregt, and the SciPy 1.0 Contributors. `scipy.stats.spearmanr` — `scipy` manual. Online doc-
1403 `scipy.stats.spearmanr` — `scipy` manual. Online doc-
umentation, SciPy, 2025. URL [https://docs.scipy.org/doc/scipy/reference/](https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.spearmanr.html)
[generated/scipy.stats.spearmanr.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.spearmanr.html). Accessed: 2025-09-21.

- 1404 Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghui Lin, Yunbo Cao, Lingpeng Kong,
1405 Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators. In Lun-Wei
1406 Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the*
1407 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9440–9450, Bangkok,
1408 Thailand, August 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024.
1409 acl-long.511. URL <https://aclanthology.org/2024.acl-long.511/>.
- 1410 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo,
1411 Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang,
1412 Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhua Chen. Mmlu-pro: A more robust and
1413 challenging multi-task language understanding benchmark. In A. Globerson, L. Mackey,
1414 D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Inform-*
1415 *ation Processing Systems*, volume 37, pp. 95266–95290. Curran Associates, Inc., 2024b.
1416 URL [https://proceedings.neurips.cc/paper_files/paper/2024/file/](https://proceedings.neurips.cc/paper_files/paper/2024/file/ad236edc564f3e3156e1b2feafb99a24-Paper-Datasets_and_Benchmarks_Track.pdf)
1417 [ad236edc564f3e3156e1b2feafb99a24-Paper-Datasets_and_Benchmarks_](https://proceedings.neurips.cc/paper_files/paper/2024/file/ad236edc564f3e3156e1b2feafb99a24-Paper-Datasets_and_Benchmarks_Track.pdf)
1418 [Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/ad236edc564f3e3156e1b2feafb99a24-Paper-Datasets_and_Benchmarks_Track.pdf).
- 1419 Jason Wei, Nguyen Karina, Hyung Won Chung, Yunxin Joy Jiao, Spencer Papay, Amelia Glaese,
1420 John Schulman, and William Fedus. Measuring short-form factuality in large language models,
1421 2024. URL <https://arxiv.org/abs/2411.04368>.
- 1422 Yuning Wu, Jiahao Mei, Ming Yan, Chenliang Li, Shaopeng Lai, Yuran Ren, Zijia Wang, Ji Zhang,
1423 Mengyue Wu, Qin Jin, and Fei Huang. Writingbench: A comprehensive benchmark for generative
1424 writing, 2025. URL <https://arxiv.org/abs/2503.05244>.
- 1426 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
1427 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
1428 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
1429 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
1430 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
1431 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
1432 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
1433 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
1434 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- 1435 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
1436 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
1437 Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on*
1438 *Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=uccHPGDlao>.
- 1440 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,
1441 and Le Hou. Instruction-following evaluation for large language models, 2023. URL <https://arxiv.org/abs/2311.07911>.

1444 A APPENDIX

- 1445
1446
1447 Table 5 lists Spearman correlation coefficients between each benchmark ranking and the ranking on
1448 each of the four task-specific categories in LMArena.
- 1449 Table 6 lists the descriptions of our adopted benchmarks.
- 1450
1451 Table 7 and Table 8 list the URLs where we collect the scores and rankings on these benchmarks.
- 1452 Table 9 lists the benchmarks we investigated but didn’t adopt. We’ll abandon a benchmark if it
1453 satisfies at least one of the descriptions below:
- 1454
- 1455 • The LLMs evaluated on the benchmark are too old that no LLMs ranking top 10 on the
1456 overall LMArena leaderboard are in its leaderboard.
 - 1457 • The benchmark has approached saturation, and no obvious distinction between benchmarks
can be seen from the benchmark.

- The benchmark is too hard and there is also no obvious distinction.
- The benchmark evaluates complex ability, which cannot be categorized into our four categories.

Table 5: Spearman correlation coefficients (ρ) between each benchmark ranking and the ranking on each of the four task-specific categories in LMArena. This table expands the schematic shown in Table 2, presenting the full set of calculated Spearman correlation coefficients ρ . Rows group the 24 benchmarks by their primary category (QA, Mathematics, Coding, Alignment). Columns represent the four LMArena rankings (Overall, Coding, Math, Instruction Following). Each cell shows the ρ value and its statistical significance ($***p < 0.001$, $**p < 0.01$, $*p < 0.05$; N.S. = Not significant). The bottom row shows the average correlation coefficient for each LMArena category across all benchmarks.

	LMArena’s Overall Ranking	LMArena’s Math Ranking	LMArena’s Coding Ranking	LMArena’s Instruction Following Ranking
QA				
SimpleQA	0.79***	0.72***	0.73***	0.80***
Facts Grounding	0.32 N.S.	0.45 N.S.	0.23 N.S.	0.27 N.S.
GPQA	0.79***	0.93***	0.73***	0.74***
MMLU-Pro	0.86***	0.86***	0.87***	0.88***
Humanity’s Last Exam	0.51 N.S.	0.83***	0.43 N.S.	0.35 N.S.
SuperGPQA	0.92***	0.94***	0.92***	0.92***
ARC-AGI-2	0.71**	0.90***	0.66**	0.71**
Math				
MGSM	0.64***	0.71***	0.65***	0.63***
MATH-500	0.72***	0.89***	0.71***	0.69***
FrontierMath (Tier 1-3)	0.78***	0.87***	0.75***	0.78***
AIME	0.58***	0.78***	0.62***	0.48**
HMMT February 2025	0.67**	0.85***	0.63**	0.66**
Coding				
HumanEval	0.74***	0.84***	0.78***	0.79***
LiveCodeBench	0.78***	0.92***	0.77**	0.66*
SWE-bench Verified	0.68*	0.54 N.S.	0.80**	0.72**
Aider polyglot	0.80***	0.84***	0.87***	0.82***
Terminal-Bench	0.80***	0.77***	0.80***	0.82***
SciCode	0.68***	0.73***	0.68***	0.72***
IOI	0.76*	0.83**	0.64*	0.78**
Instruction Following				
IFEval	0.34 N.S.	0.39 N.S.	0.31 N.S.	0.45*
IFBench	0.58***	0.70***	0.59***	0.61***
ArenaHard	0.73***	0.94***	0.80***	0.73***
WritingBench	0.71***	0.78***	0.79***	0.72***
Creative Writing v3	0.80***	0.70***	0.79***	0.84***
Average	0.69	0.78	0.69	0.69

Table 6: Detailed Descriptions of Benchmark Datasets and Their Evaluation Focus

Benchmark	Description
SimpleQA	Short, simple questions answering in broad domains.
FACTS Grounding	Responses generation based on provided documents, testing factuality especially in critical domains.
GPQA	Ultra-hard multiple-choice problems testing beyond-human knowledge.
MMLU-Pro	Reasoning, comprehension and memorizing abilities evaluation with 12K+ complex questions in 14 subjects.
Humanity’s Last Exam	2,500 Advanced reasoning test using 2500 extremely difficult questions over more than ten subjects.
SuperGPQA	Graduate-level questions reasoning covering niche professional fields.
ARC-AGI-2	Evaluation of cognitive flexibility using grid-based puzzles, which requires ability of abstract reasoning, pattern cognition and generalization.
MGSM	Multilingual evaluation using translated GSM8K problems.
MATH-500	500-problem collection from MATH, spanning diverse mathematical knowledge.
FrontierMath	Challenging mathematical problems covering most major branches of modern mathematics.
AIME	Solving of high-difficulty competition problems covering algebra, geometry, number theory, and combinatorics.
HMMT February 2025	Problems from an undergraduate-level math tournament.
HumanEval	A set of 164 hand-written programming problems for measuring functional correctness for synthesizing programs from docstrings.
LiveCodeBench	Code generation, repair, and comprehension on 300+ recent problems.
SWE-bench Verified	Benchmark of abilities to solve real-world software issues by evaluating agents’ ability to generate a patch for given GitHub repository and issue description.
Aider Polyglot	Low-solve-rate coding problems from Exercism in six languages.
Terminal-Bench	Tasks for evaluating agents’ performance in using terminals.
SciCode	Challenging problems from 16 diverse natural science fields to evaluate LMs’ capabilities of code generation for the solution of real scientific research problems.
IOI	Problems from International Olympiad in Informatics
IFEval	Instruction-following capability testing with 500+ diverse prompts.
IFBench	Benchmark to evaluate precise instruction following generalization on 58 new, diverse, and challenging verifiable out-of-domain constraints.
ArenaHard	High-difficulty real-user queries testing ability of interacting.
WritingBench	Writing ability evaluation across 100+ article prompts.
Creative Writing v3	Expansion of Creative Writing v2, using a hybrid rubric and Elo scoring system to enhance discrimination.

Table 7: Details of the Sources of Accuracy Data on Benchmark Datasets

Benchmark	Details
SimpleQA	Scores are obtained from https://www.kaggle.com/benchmarks/openai/simpleqa , and the last updated date of which is 2025/09/03.
FACTS Grounding	Scores are obtained from the column marked with caption “Score” in https://www.kaggle.com/benchmarks/google/facts-grounding , and the last updated date of which is 2025/08/20.
GPQA	Results are obtained from the “Overall” task type in https://www.vals.ai/benchmarks/gpqa-09-08-2025 .
MMLU-Pro	Results are obtained from the “Overall” column in https://huggingface.co/spaces/TIGER-Lab/MMLU-Pro .
Humanity’s Last Exam	Accuracy datas were fetched from https://agi.safe.ai/ , last updated on 2025/04/03.
SuperGPQA	Results are obtained from the “Overall” column in https://supergpqa.github.io/ .
ARC-AGI-2	Results were obtained from column marked “ARC-AGI-2” in https://arcprize.org/leaderboard .
MGSM	Results are obtained from the “Overall” task type in https://www.vals.ai/benchmarks/mgsm-2025-09-08 .
MATH-500	Results are obtained from the “Overall” task type in https://www.vals.ai/benchmarks/math500-05-30-2025 .
FrontierMath	Accuracy data are obtained from tab “Tier 1-3” in https://epoch.ai/frontiermath . Note that manual switching to the tab is required.
AIME	Results are obtained from the “Overall” task type in https://www.vals.ai/benchmarks/aime-2025-09-08 .
HMMT Feburary 2025	Accuracy data are obtained from tab “HMMT Feb 2025” in https://matharena.ai/?comp=hmmt-hmmt_feb_2025 .
HumanEval	Results are obtained from the “HUMANEVAL” column in https://artificialanalysis.ai/leaderboards/models , note that the column will be visible only after clicking the “Expand Columns” button to the right of the page.
LiveCodeBench	Results are obtained from the “PASS@1” column in https://livecodebench.github.io/leaderboard.html , note that the time window should be manually adjusted to 5/1/2023-5/1/2025.
SWE-bench Verified	Data were obtained from tab “Bash Only” in https://www.swebench.com/ .
Aider Polyglot	Results are obtained from the “Percent correct” column in https://aider.chat/docs/leaderboards/ .
Terminal-Bench	Results are obtained from the “Overall” task type in https://www.vals.ai/benchmarks/terminal-bench-2025-09-18 .
SciCode	Results are obtained from the “SCICODE” column in https://artificialanalysis.ai/leaderboards/models , note that the column will be visible only after clicking the “Expand Columns” button to the right of the page.
IOI	Results are obtained from the “Overall” task type in https://www.vals.ai/benchmarks/ioi-09-09-2025 .

Table 8: Details of the Sources of Accuracy Data on Benchmark Datasets (Continued)

Benchmark	Details
IFEval	Results are obtained from the “IFEval - IFEval Strict Acc” column in https://crfm.stanford.edu/helm/capabilities/latest/#/leaderboard .
IFBench	Results are obtained from the “IFBENCH” column in https://artificialanalysis.ai/leaderboards/models , note that the column will be visible only after clicking the “Expand Columns” button to the right of the page.
ArenaHard	Results are obtained from the “Official Configuration” sector in README.md of https://github.com/lmarena/arena-hard-auto?tab=readme-ov-file#leaderboard .
WritingBench	Results are obtained from the “Overall” column in https://huggingface.co/spaces/WritingBench/WritingBench .
Creative Writing v3	Results are obtained from column “Rubric Score” in https://eqbench.com/creative_writing.html .

Table 9: Detailed Descriptions of Benchmark Datasets that were Examined but Not Used

Benchmark	Description
MMLU	Multi-subject, broad-difficulty question answering, testing problem-solving depth and breadth.
BBH	Hard problems reasoning, enforcing CoT.
MuSR	Multistep soft reasoning tasks specified in a natural language narrative.
ZebraLogic	Non-monotonic logical reasoning under constraint satisfaction problems.
τ^2 -Bench	Evaluating conversational agents in a dual-control environment where both the agent and the user can operate tools, requiring coordination and communication between the agent and the user to succeed.
MATH	12500 challenging competition math problems with step-by-step solution.
GSM8K	Elementary-level math word problems assessing basic reasoning.
MathBench	Progressive math from arithmetic to university-level, testing theoretical and applied skills.
MultiPL-E	Code generation over broad programming languages using problems in HumanEval and MBPP.
EvalPlus	Rigorous correctness testing using HumanEval and MBPP advanced by adversarial mutations.
CRUXEval	Code semantic understanding via output prediction and backward reasoning.
InfiBench	Code comprehension and correlated natural language QA in real context.
EvoEval	Evolution-augmented problems solving reducing overfitting.
MHPP	210 Python problems with 7 core challenges testing NLP and edge-case handling.
BigCodeBench	Real-world challenging code completion and generation tasks.
AlpacaEval 2.0	Advanced instruction-following benchmark with anti-gaming measures.
Creative Writing v2	Text generation quality assessment through diverse prompts.