

LLM-Powered Benchmark Factory: Reliable, Generic, and Efficient

Anonymous Authors¹

Abstract

The rapid advancement of large language models (LLMs) has led to a surge in both model supply and application demands. To facilitate effective matching between them, reliable, generic and efficient benchmark generators are widely needed. However, human annotators are constrained by inefficiency, and current LLM benchmark generators not only lack generalizability but also struggle with limited reliability, as they lack a comprehensive evaluation framework for validation and optimization. To fill this gap, we first propose an automated and unbiased evaluation framework, structured around four dimensions and ten criteria. Under this framework, we carefully analyze the advantages and weaknesses of directly prompting LLMs as generic benchmark generators. To enhance the reliability, we introduce a series of methods to address the identified weaknesses and integrate them as BENCHMAKER. Experiments across multiple LLMs and tasks confirm that BENCHMAKER achieves superior or comparable performance to human-annotated benchmarks on all metrics, highlighting its generalizability and reliability. More importantly, it delivers highly consistent evaluation results across 12 LLMs (0.967 Pearson correlation against MMLU-Pro), while taking only \$0.005 and 0.38 minutes per sample.

1. Introduction

With the ongoing scaling up of large language models (LLMs) in multiple dimensions over the past few years, two key trends have emerged (Figure 1): (1) The LLM release process has accelerated and now exceeds 30k per season; (2) The growth in LLM capabilities has spurred application demand, reflected in over 50M downloads of open-source models per season. Serving as a bridge be-

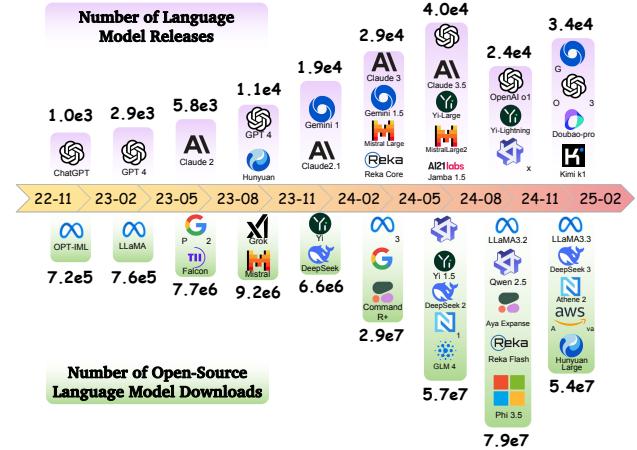


Figure 1. The trends of LLMs released and open-source LLMs downloads per season since the debut of ChatGPT. We obtain the data via the Huggingface API. See details in Appendix C.

tween massive LLM supply and various application needs, the demand for customized benchmarks is rapidly growing, helping downstream tasks identify the most suitable LLM.

However, current benchmark construction processes largely rely on human-provided signals (Chang et al., 2024; Wang et al., 2024b), leading to long cycles and high costs. To this end, efficient LLM-driven methods have recently been explored. Unfortunately, they generally rely on the existence of seed benchmarks for data augmentation (Zhu et al., 2024b; Wu et al., 2024; Li et al., 2024a; Maheshwari et al., 2024) and task specific designs (Zhu et al., 2024a; Lei et al., 2023), lacking generalization across tasks and domains. Meanwhile, the current absence of a comprehensive evaluation framework hinders the assessment and optimization of benchmark generators, weakening our confidence in their reliability for real applications. Hence, an automatic and comprehensive evaluation framework and a generic and reliable benchmark generator that can handle any assessment demands and efficiently generate high-quality samples are urgently needed.

To this end, we first construct an automatic evaluation framework with ten criteria for benchmark generators. Notably, we utilize causal learning (Kaddour et al., 2022) techniques to identify and remove biases of LLM-as-a-judge (Liu et al., 2023) across various criteria, ensuring the reliability of the

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review . Do not distribute.

framework. On this basis, we examine the strengths and weaknesses of directly prompting LLM as generic benchmark generator through this evaluation framework. The results reveal that the generated benchmark exhibits limited lexical and semantic diversity, poor controllability over difficulty, and low sample faithfulness, while showing advantages in high task alignment and knowledge diversity. Bearing this in mind, we develop a generic benchmark generator BENCHMAKER by integrating existing techniques with newly designed approaches to address the identified issues. Specifically, BENCHMAKER: strengthens sample faithfulness using stepwise self-correction generation and conflict guided contrastive discrimination; extends difficulty boundary with difficulty strategy guidance and difficulty diffusion mechanism; enhances diversity through AttrPrompt (Yu et al., 2023) and in-batch redundancy filtering. We also discuss some unsuccessful attempts in Appendix B to provide more insights for future research.

We conduct comprehensive experiments to validate BENCHMAKER under the proposed evaluation framework. Compared to high-quality human-annotated benchmarks, the benchmarks generated by BENCHMAKER exhibit superior task alignment, better difficulty controllability, more challenging question difficulties, and comparable sample faithfulness and diversity. More importantly, they yield highly consistent evaluation results across 12 LLMs (0.967 Pearson correlation with MMLU-Pro), with BENCHMAKER taking only \$0.005 and 0.38 minutes per sample. We further perform detailed experiments to validate the outstanding generalization and robustness across tasks and LLMs, and the effectiveness of each component of BENCHMAKER. Finally, we derive a formula for evaluating the confidence of benchmarking results under conditions where faithfulness cannot be fully satisfied, further enhancing the practicality and reliability of BENCHMAKER.

2. Backgrounds

In this section, we first review the latest developments in data synthesis §2.1 and then discuss the potential values of developing a generic benchmark generator §2.2.

2.1. Synthetic Data Generation

The growth of language model abilities has led to widespread research on LLM-driven data synthesis, which demonstrates much better quality and controllability over traditional approaches (Wang et al., 2024a; Long et al., 2024). Centering around the construction of data flywheel (LLM-driven evolution) (Luo et al., 2024a; Tao et al., 2024), training data synthesis has garnered much attention in fields like mathematics (Yu et al., 2024), science (Li et al., 2024b), and code (Luo et al., 2024b), continuously pushing LLMs’ capability boundaries. Unlike the training data synthesis

aimed at optimizing model performance, the goal of benchmark synthesis is to accurately evaluate models on specific task, presenting greater challenges in both measurement and implementation (Chang et al., 2024). In terms of measurement, recent studies (Zhu et al., 2024a; Maheshwari et al., 2024; Li et al., 2024a) generally focus on specific criteria, without establishing a comprehensive evaluation system for benchmark generators. In terms of implementation, current benchmark generators (Perez et al., 2023; Wu et al., 2024; Zhu et al., 2024b; Lei et al., 2023) are constrained by their dependence on existing benchmarks and task specific designs, preventing them from being generic. We construct a comprehensive evaluation framework and develop generic and reliable BENCHMAKER method to fill this gap.

2.2. Potential Applicable Scenarios of BENCHMAKER

Given arbitrary assessment demands X as the sole input, a generic benchmark generator (BENCHMAKER) \mathcal{G} is expected to generate a well-aligned high-quality benchmark \mathcal{D} . On this basis, we summarize its applicable scenarios as follows: (1) Complementing existing benchmarks for tailored assessment demands; (2) Acting as a dynamic benchmark generator to alleviate data contamination issues (Balloccu et al., 2024); (3) Serving as a difficulty controllable benchmark generator to mitigate the benchmark saturation problem (Glazer et al., 2024); (4) Functioning as a versatile training data generator. Therefore, building BENCHMAKER holds significant importance for both scientific research and practical applications within the NLP community.

3. Benchmarking Benchmark Generator

While training data synthesis focuses on faithfulness, diversity and the final performance of the trained models (Yu et al., 2023; Long et al., 2024), the evaluation of synthetic benchmark should be more comprehensive to ensure the reliability of its benchmarking results. Thus, we carefully establish an evaluation framework for benchmark generator with ten criteria, as illustrated in Table 1.

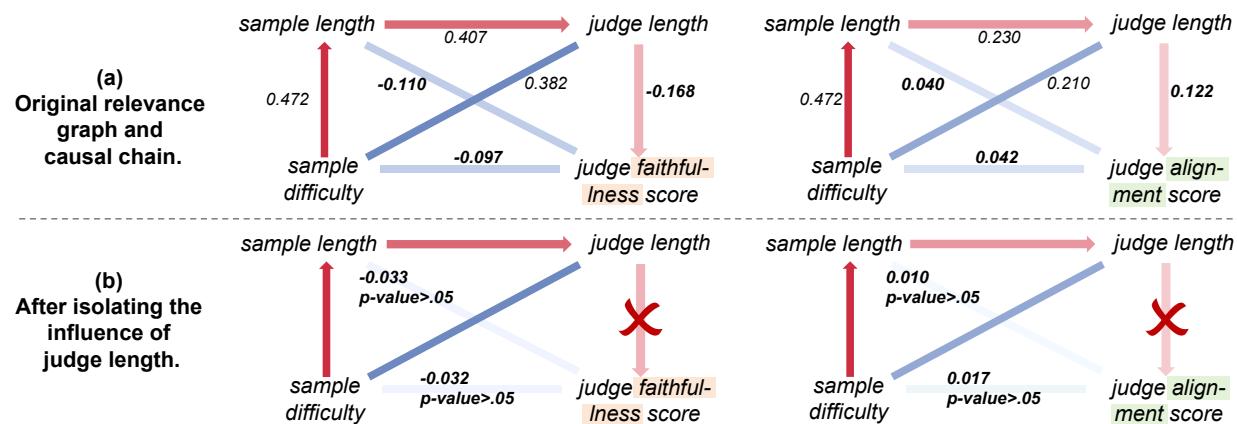
3.1. Credibility

Two key criteria for ensuring the credibility of a benchmark are **faithfulness** and **alignment**. Faithfulness indicates that the generated sample be free of ambiguity with a correct answer. Alignment requires the generated samples to strictly adhere to the specified assessment demands X , especially in abilities to be assessed. For these criteria, previous approaches rely on human evaluation (Wu et al., 2024; Zhu et al., 2024b) or LLM-as-a-judge (Zheng et al., 2023). However, the former lacks automation, and the latter is susceptible to biases (Thakur et al., 2024).

To this end, we seek to detect and mitigate any biases of

110
111 **Table 1.** Criteria taxonomy and definition of the proposed evaluation framework for the benchmark generator. Criteria marked with *
112 indicate optimization objectives that are distinctive to benchmark synthesis compared to training data synthesis.

Taxonomy	Criterion	Definition
Credibility	faithfulness	The sample is well-defined, with the ground truth answer being correct.
	alignment *	The abilities evaluated by the sample align well with the given assessment demands.
Diversity	lexical	The samples exhibit sufficient lexical diversity.
	semantic	The samples exhibit sufficient semantic richness.
Difficulty	knowledge *	The knowledge and skills assessed by different samples should not be redundant.
	controllability *	The samples have correct difficulty labels to form subsets with varying difficulties.
Benchmark-Level	boundary *	The hardest subset is difficult enough to explore the boundaries of advanced models.
	effectiveness *	The benchmarking results align with human benchmark under the same assessment demands.
Benchmark-Level	robustness *	The benchmarking results of generated benchmarks under similar assessment demands align.
	efficiency	The time and cost of generating a benchmark are low enough.



133 **Figure 2.** Pearson correlations among key factors of benchmark evaluation and LLM (Qwen-Plus) judge scores (faithfulness and alignment).
134 The most relevant path of each subject is highlighted in red to show the possible causal chain.
135
136
137
138

139 LLM-as-a-judge that may exist within the framework. We
140 choose Qwen-Plus (Yang et al., 2024) as the judge with
141 scoring range as [0, 1] (See prompt in Appendix I). Experi-
142 ments are conducted on the high-quality MATH benchmark
143 (Hendrycks et al., 2021b), for which we assign score 1 to
144 both faithfulness and alignment for every sample. Ideally,
145 the scores assigned by the judge should not exhibit any
146 consistency with specific factors. However, as shown in Fig-
147 ure 2-(a), both faithfulness and alignment are significantly
148 correlated ($p\text{-value} < 0.05$) with sample difficulty, sample
149 length, and the length of the judge’s rationale. For each
150 factor, we highlight its weightiest path in red, revealing a
151 clear causal chain: harder questions lead to longer samples,
152 requiring judges to conduct lengthier analyses. For faith-
153 fulness, longer analyses increases the likelihood of judge
154 errors, resulting in lower faithfulness ratings. While for
155 alignment, longer analyses increases the probability of task-
156 relevant words appearing and results in higher alignment
157 ratings. To validate the above hypothesis, we control the
158 judge length and respectively calculate the partial corre-
159 lations (Vallat, 2018) of sample difficulty and sample length
160 with faithfulness and alignment. As shown in Figure 2-(b),
161 apart from the path from sample length to judge length,
162 the path from sample difficulty to judge length is also
163 highlighted in red, indicating that sample difficulty has
164 a significant influence on judge length even after controlling
165 for judge length.

166 after isolating the influence of judge length, the effects of
167 other factors are no longer significant ($p\text{-value} > 0.05$). Similar
168 conclusions also hold true when GPT-4o mini (Hurst
169 et al., 2024) serves as the judge (Figure 7 in Appendix).

170 Based on the analysis above, the potential biases of the
171 LLM judge in this scenario are all mediated by judgment
172 length. Therefore, for benchmark generators $\mathcal{G}_{1:|\mathcal{G}|}$ under
173 evaluation, we derive their unbiased judge results with a
174 Multiple Regression model. Specifically, we set the judge
175 score as the dependent variable, the generator categories as
176 dummy variables, and judge length as the covariate:

$$f(i) = \beta_i + \beta_{len} \cdot \text{judge.length} + \epsilon \quad (1)$$

177 where $f(i)$ denotes the average judge score of \mathcal{G}_i and β_i
178 reflects the debiased score of \mathcal{G}_i , which we select as our
179 metrics for faithfulness and alignment.

3.2. Diversity

180 With credibility ensured, the diversity of the benchmark
181 determines the extent to which evaluation results can reflect
182 the true model capability across the assessed domain. Apart
183 from the LLM judge, we also evaluate the diversity of the
184 generated benchmarks using the NLP diversity metric (Liu
185 et al., 2021) and the BLEU score (Papineni et al., 2002).
186

165 from the widely tested lexical and semantic diversity, our
 166 framework also examines the knowledge diversity to make
 167 the evaluation more comprehensive.
 168

169 **Lexical Diversity** reflects vocabulary richness in benchmarks.
 170 Traditional metrics like vocabulary size and self-
 171 BLEU (Zhu et al., 2018) used in Wu et al. (2024) and (Yu
 172 et al., 2023) are biased by sample length (Guo & Vosoughi,
 173 2023). We use unbiased word frequency entropy (Montahaei
 174 et al., 2019) as the metric to evaluate lexical diversity.
 175

176 **Semantic Diversity** quantifies a benchmark’s semantic
 177 comprehensiveness. We calculate the average Euclidean
 178 distance between semantic embeddings of samples as the
 179 metric. Specifically, we use powerful text-embedding-ada-
 180 002 (OpenAI) as the embedding model.
 181

182 **Knowledge Diversity** evaluates whether the samples eval-
 183 uate different sub-abilities within the assessment demands.
 184 When samples test the same sub-ability, the model is likely
 185 to exhibit similar correctness patterns. Therefore, we use
 186 the correctness of a set of models (denoted as $\mathcal{M}_{1:|\mathcal{M}|}$, see
 187 Appendix E for detailed list) to represent the knowledge
 188 embedding for each sample. If the embeddings of two sam-
 189 ples are highly similar, it reflects a strong alignment in the
 190 sub-abilities they assess. The average pairwise Hamming
 191 distance (Hamming, 1950), suitable for discrete embeddings,
 192 is employed as the metric for this criterion.
 193

194 3.3. Difficulty

195 When diversity meets requirements, we should further con-
 196 sider the difficulty attribute, which is particularly significant
 197 in an era of increasingly divergent model capabilities.
 198

199 **Difficulty Controllability** refers to assigning differenti-
 200 ated difficulty labels to the samples (e.g., MATH (Hendrycks
 201 et al., 2021b)). These labels enable the benchmark to be
 202 divided into subsets for more targeted evaluation of models
 203 with varying capabilities. For each sample, we use the aver-
 204 age error rate of $\mathcal{M}_{1:|\mathcal{M}|}$ as the ground truth for difficulty
 205 label. Based on this, we compute the Spearman correlation
 206 between the difficulty labels provided by the benchmark and
 207 the ground truth as the metric.
 208

209 **Difficulty Boundary** denotes the difficulty of the hardest
 210 subset of a benchmark. With the growing strength of LLMs,
 211 their performance on simpler benchmarks has reached sat-
 212 uration (Hendrycks et al., 2021a), making it difficult to
 213 differentiate their capabilities. Consequently, more chal-
 214 lenging benchmarks (Wang et al., 2024b) are continuously
 215 introduced to evaluate the latest LLMs. Thus, we propose
 216 assessing the average error rate of $\mathcal{M}_{1:|\mathcal{M}|}$ on the hardest
 217 subset of benchmark to measure its difficulty boundary.
 218

3.4. Benchmark-Level

Lastly, we introduce high-level metrics for assessing bench-
 mark generators.

Effectiveness. While the earlier criteria assess benchmark
 quality from various aspects, a unified metric is required
 to measure benchmark effectiveness. Taking high-quality
 human-annotated benchmark as the ground truth, we exam-
 ine whether generated benchmark under identical assess-
 ment demands can deliver equivalent evaluation results. To
 this end, we calculate the accuracy of $\mathcal{M}_{1:|\mathcal{M}|}$ on both gener-
 ated and human benchmarks and use the Pearson correlation
 between them as the effectiveness metric.

Robustness. Under similar inputs, a robust system should
 produce comparable outputs. Similarly, we expect a robust
 benchmark generator to produce benchmarks with equiv-
 alent evaluation efficacy for similar assessment demands.
 Thus, we calculate the accuracy of $\mathcal{M}_{1:|\mathcal{M}|}$ on benchmarks
 generated under similar assessment demands (the original
 and that rewritten by GPT-4o) and calculate the Pearson
 correlation between them as the robustness metric.

Efficiency. High-quality human-annotated benchmarks
 are constrained due to inefficiencies in their construction.
 We evaluate the efficiency of a benchmark generator by
 measuring the time and monetary costs associated with gen-
 erating benchmarks of a certain size.

By establishing this comprehensive evaluation framework,
 the strengths and weaknesses of benchmark generators can
 be thoroughly assessed, and the reliability of the proposed
 method can be validated.

4. Development of BenchMaker

In this section, we first discuss the primary sample format
 we studied in §4.1. Afterwards, since previous studies have
 yet to realize generic benchmark generators (with assess-
 ment demands X as the sole input), we analyze the pros
 and cons of directly prompting the LLM as such generator
 in §4.2. Building on the experimental results, we refine its
 weaknesses in the following sections, leading to the devel-
 opment of BENCHMAKER.

4.1. Sample Format Selection

Following previous studies (Li et al., 2024a; Zhu et al.,
 2024b), we have chosen multiple-choice questions (MCQs)
 as the primary sample format for benchmark generation
 based on the following reasons: (1) Versatility: MCQ serves
 as a universal format for evaluating most capabilities; (2)
 Accuracy: Misjudgment can be effectively prevented caused
 by variations in output formats (Tam et al., 2024); (3) Effi-

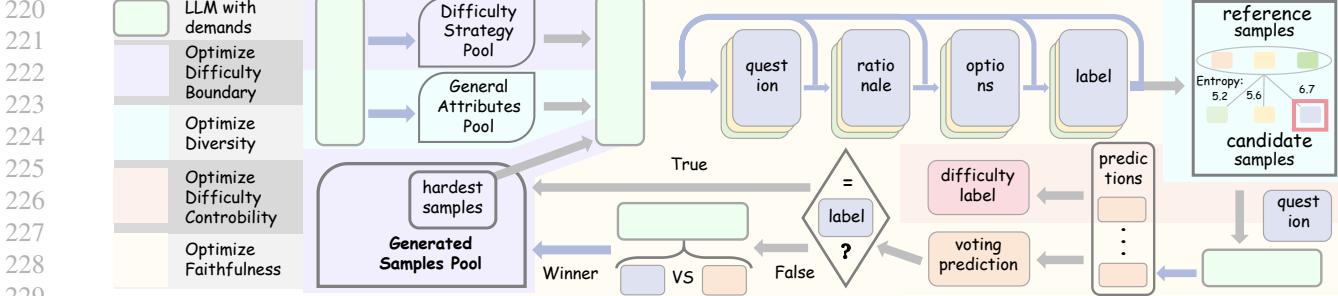


Figure 3. Overview of BENCHMAKER.

ciency: MCQs do not depend on external modules such as LLM-as-a-Judge, ensuring a streamlined evaluation process; (4) Transformability: Each generated sample includes a rationale, enabling easy conversion into other formats, such as the text generation format presented in the Appendix G.

4.2. Pros and Cons of Directly Prompting

We choose MATH (Hendrycks et al., 2021b), MMLU-Pro (Wang et al., 2024b) and HellaSwag (Zellers et al., 2019) as high-quality benchmarks \mathcal{D}_{human} for comparison. We adopt the prompt in Appendix I to guide \mathcal{M} : GPT-4o mini in generating credible and diverse samples $s_{1:|\mathcal{D}_{human}|}$:

$$s_i = \{q_i, r_i, o_i, a_i\} = \mathcal{M}(\text{prompt}_{base}, l, X) \quad (2)$$

where q_i, r_i, o_i, a_i denote question, rationale, options, and label, respectively. We proportionally adjust the difficulty level l from 1 to 10 in the prompt (see descriptions in Appendix D), and select samples with top 20% difficulty level to form the hardest subset. The assessment demands X are shown in Appendix K. As shown in Table 2, compared to \mathcal{D}_{human} (red line), directly prompting LLM as generic benchmark generator (yellow line) demonstrates poorer faithfulness, lower lexical and semantic diversity, weaker difficulty controllability, and less challenging subset. Meanwhile, we also observe its advantages in better alignment¹, greater knowledge diversity, and improved efficiency.

4.3. Faithfulness Optimization

To enhance faithfulness, previous studies have explored methods such as self-correction (Wang et al., 2023b; Ji et al., 2023) and the use of external tools (Li et al., 2024c; Lewis et al., 2020). As self-correction offers greater versatility, we propose the following two BenchMaker-compatible techniques to optimize faithfulness.

¹We set the debiased LLM-as-a-judge score of the human benchmark to 1, adjusting scores of generated benchmarks accordingly, which may result in scores exceeding 1.

Stepwise Self-correction. Since errors might occur at any step during the generation of $\{q_i, r_i, o_i, a_i\}$, we instruct the model to validate the content at each step. If an error is detected, the model will return to the beginning. Compared to full-sample self-checking, step-wise critique boosts error detection with less decoding cost (See Appendix B).

Conflict Guided Contrastive Discrimination. Huang et al. (2024) finds that LLMs struggle to correctly judge their prior answers on challenging questions. Therefore, we extend Stepwise Self-correction by having the LLM not only act as a judge but also as a test-taker to identify potential errors. Let the LLM answers $q_i T$ times to attain $\bar{a}_i^{1:T}$, we get the self-consistency (Wang et al., 2023a) result \hat{a}_i through majority voting. If $\hat{a}_i \neq a_i$, the conflict suggests differing r_i and \hat{r}_i . As Zheng et al. (2023) finds that comparison-based judges are more accurate than item-wise judges, we have the LLM conduct a contrastive discrimination between r_i and \hat{r}_i to determine the final rationale and label for s_i .

4.4. Difficulty Optimization

Difficulty Controllability. From §4.2, we know that the LLM’s ability to control the difficulty of generated samples is limited. In particular, for the language understanding task (MMLU-Pro), the Spearman correlation between the actual and expected difficulty of the samples is only 0.021. To further explore this, we examine LLM’s difficulty perception by asking it to score the difficulty label of the generated samples. However, the correlation only increases to 0.089, suggesting that while LLM has some capacity to perceive difficulty, it is still weak. We then switch the role of LLM and assess the difficulty from the perspective of test-taker:

$$\beta_i = \frac{1}{T} \sum_{j=1}^T \mathbf{1}_{\bar{a}_i^j \neq a_i} \quad (3)$$

By taking the inconsistency between $\bar{a}_i^{1:T}$ and a_i as difficulty label, the correlation increases to 0.415, suggesting that β is a reliable metric for difficulty controllability.

275 **Difficulty Diffusion Mechanism.** Given that the LLM
 276 has a certain level of difficulty perception, we iteratively
 277 select the more challenging samples according to β from
 278 the generated ones as difficulty references, and instruct the
 279 LLM to generate a more difficult sample. This allows the
 280 sample difficulty to rise continuously through diffusion. The
 281 detailed algorithm is described in Appendix F.
 282

283 **Difficulty Strategy Guidance.** We further consider pro-
 284 viding the LLM with task-specific difficulty-control strate-
 285 gies. Specifically, we first require the LLM to give varying
 286 strategies for generating samples of specific difficulty levels
 287 based on the given X (see examples in Appendix J). For
 288 example, difficult samples those assessing reasoning ability
 289 generally require more reasoning steps. With the Difficulty
 290 Diffusion Mechanism, we progressively introduce more dif-
 291 ficult sample generation strategies to the LLM to further
 292 extend the difficulty boundary.
 293

294 4.5. Diversity Optimization

295 The optimization of synthetic data diversity has been widely
 296 studied (Wang et al., 2024a). We conduct extensive tests and
 297 select the most generic and effective **AttrPrompt** (Yu et al.,
 298 2023) technique for BENCHMAKER. AttrPrompt explicitly
 299 enhances the lexical and semantic diversity of benchmarks
 300 by randomly assigning pre-generated (attribute, value) pairs
 301 as part of the input for each sample. Furthermore, we notice
 302 that the introduction of treating the generated samples as
 303 difficulty references might cause sample homogeneity. To
 304 mitigate this, we propose an **In-batch Diversity Boosting**
 305 method, where LLM generates L (We set L as 5 for our
 306 default setting) candidate samples and selects the one with
 307 the greatest word frequency entropy difference from the
 308 input reference samples.
 309

310 5. Experiments and Analyses

311 We conduct comprehensive experiments to validate BENCH-
 312 MAKER under the proposed framework in this section.
 313

314 **Settings.** We select the widely used human-annotated
 315 MATH² (Hendrycks et al., 2021b) (mathematical reasoning),
 316 MMLU-Pro (multi-task language understanding) (Wang
 317 et al., 2024b) and HellaSwag (commonsense reasoning)
 318 (Zellers et al., 2019) as high-quality baseline benchmarks.
 319 For the 7 subsets of MATH, the 13 subsets of MMLU-Pro³
 320 and HellaSwag, we write simple assessment demands re-
 321 spectively (see details in Appendix K) as inputs for the
 322 benchmark generator. For each demand, we generate 500
 323 samples and randomly downsample the human-annotated
 324 benchmark to match the number of generated samples for
 325

²Converted into a MCQ format, see details in Appendix G.

³Excluding the type ‘other’.

fair comparison. Each experiment is repeated three times, and the average results are reported. We use GPT-4o mini (Hurst et al., 2024) as the default generator and also explore the performance of GPT-4o and Claude 3.5 Haiku (Anthropic). The decoding temperature is set to 1. To mitigate the self-enhancement bias (Zheng et al., 2023) associated with LLM-as-a-judge, we substitute the generators with Qwen-Plus (Yang et al., 2024) as the judge.

5.1. Comparison with Human-annotated Benchmark

As shown in Table 2, overall, BENCHMAKER achieves comparable performance to human-annotated benchmarks in terms of faithfulness and lexical&semantic diversity. Meanwhile, BENCHMAKER outperforms them in all other metrics, especially in alignment, knowledge diversity, difficulty controllability and efficiency. The exceptional results achieved in these metrics comprehensively validate the reliability of the generated samples by BENCHMAKER.

Effectiveness. The primary goal of benchmarking is to assign accurate scores to models under evaluation, facilitating capability differentiation. The benchmarking results of BENCHMAKER align closely with human-annotated benchmarks, with an average of 0.953 linear correlation (Pearson) and a remarkable 0.966 for rank-order correlation (Spearman), highlighting its outstanding effectiveness.

Robustness. Under evaluation demands where semantic equivalence is maintained but linguistic styles vary, the benchmarks exhibit nearly identical assessment efficacy, with an average Pearson correlation of 0.984. This demonstrates the robustness of BENCHMAKER to diverse inputs and ensures that users with different linguistic preferences can obtain consistent evaluation results.

Efficiency. The primary limitation of human-annotated benchmarks lies in their low construction efficiency. However, BENCHMAKER can generate a sample at an average cost of \$0.005 within 0.40 minutes. Furthermore, its efficiency is expected to continuously improve with the development of technology and hardware.

Generalizability. Experimental results demonstrate that BENCHMAKER exhibits strong generalization across different task types and generators. Notably, a more powerful model does not necessarily yield superior performance across all metrics. Compared to GPT-4o, GPT-4o mini proves to be a more cost-effective benchmark generator.

5.2. Ablation Studies

We validate the effectiveness of different techniques by sequentially integrating them to the Direct Prompt baseline on the MATH benchmark, as shown in Table 2.

330
 331 *Table 2.* Overall experimental results under the proposed evaluation framework. For each setting, we run three times and report the average
 332 results. We take GPT-4o mini as default generator. Values in bold denote the best results between BENCHMAKER and Human Benchmark.

333 334 335 336 Methods	Faithful	Alignment	Lexical	Semantic	Knowledge	Control	Boundary	Effective	Robust	Efficiency
	Unbias Score↑	Unbias Score↑	Entropy↑	Euclidean Distance↑	Hamming Distance↑	Spearman↑	Error Rate↑	Pearson↑	Pearson↑	\$/item, min/item↓
MATH (Hendrycks et al., 2021b)										
Human Benchmark	1.000	1.000	8.054	0.665	0.349	0.143	0.752	-	-	high
Direct Prompt	0.665	1.166	7.091	0.618	0.365	0.109	0.635	0.687	0.991	0.002, 0.17
+AttrPrompt	0.611	1.138	8.265	0.675	0.360	0.124	0.659	0.759	0.983	0.002, 0.19
+InBatchDivBoost	0.623	1.142	8.652	0.677	0.366	0.115	0.628	0.778	0.985	0.003, 0.20
+StepSelfCorrect	0.924	1.152	8.674	0.675	0.369	0.162	0.557	0.803	0.979	0.003, 0.23
+ConflictConDisc	1.019	1.151	8.668	0.678	0.357	0.175	0.515	0.838	0.992	0.004, 0.35
+DiffControl	1.019	1.151	8.668	0.678	0.357	0.403	0.515	0.838	0.992	0.004, 0.35
+DiffDiffusion	0.994	1.166	8.705	0.680	0.387	0.451	0.683	0.882	0.990	0.005, 0.39
BenchMaker	0.930	1.200	8.976	0.681	0.403	0.434	0.768	0.935	0.986	0.005, 0.42
BenchMaker _{4o}	0.918	1.223	8.835	0.675	0.385	0.432	0.779	0.941	0.988	0.084, 1.12
BenchMaker _{haiku}	0.902	1.116	8.878	0.676	0.410	0.401	0.775	0.912	0.979	0.026, 0.57
MMLU-Pro (Wang et al., 2024b)										
Human Benchmark	1.000	1.000	10.404	0.731	0.307	0.000	0.751	-	-	high
Direct Prompt	0.894	1.218	9.608	0.726	0.391	0.021	0.587	0.850	0.989	0.002, 0.16
BenchMaker	1.020	1.245	10.166	0.728	0.395	0.477	0.759	0.967	0.982	0.005, 0.38
HellaSwag (Zellers et al., 2019)										
Human Benchmark	1.000	1.000	9.167	0.655	0.384	0.000	0.569	-	-	high
Direct Prompt	0.862	1.107	8.165	0.660	0.396	0.047	0.626	0.821	0.979	0.002, 0.17
BenchMaker	1.032	1.130	9.052	0.663	0.421	0.439	0.708	0.958	0.984	0.005, 0.40

358 **Diversity.** Compared to Direct Prompt, both AttrPrompt
 359 and In-batch Diversity Boosting effectively enhance lexical
 360 and semantic diversity. Noticeably, we observe that knowl-
 361 edge diversity remains unchanged, indicating that surface-
 362 level diversification does not necessarily equate to a broader
 363 assessment of knowledge and skills. Meanwhile, the di-
 364 versity improvement leads to a slight drop in faithfulness,
 365 possibly because of the attributes constraints.

366 **Faithfulness.** After applying Stepwise Self-correction and
 367 Conflict Guided Contrastive Discrimination, we observe a
 368 sustained and significant improvement in faithfulness. At
 369 the same time, we notice a reduction in the difficulty of the
 370 hardest subset, with the error rate decreasing from 0.659 to
 371 0.557. We hypothesize that this may be due to the high error
 372 rate in labels when faithfulness is not ensured, which leads
 373 to an underestimation of model performance. Consequently,
 374 once the labels are corrected, the accuracy can better reflect
 375 the actual difficulty of the benchmark.

376 **Difficulty Controllability.** By treating the generator as
 377 the test-taker and using its error rate as the difficulty la-
 378 bel, we achieve more precise control over sample difficulty
 379 (Spearman correlation of 0.403). Considering the previously
 380 observed weak difficulty perception of LLMs, we hypothe-
 381 size that this improvement stems from the role shift, which

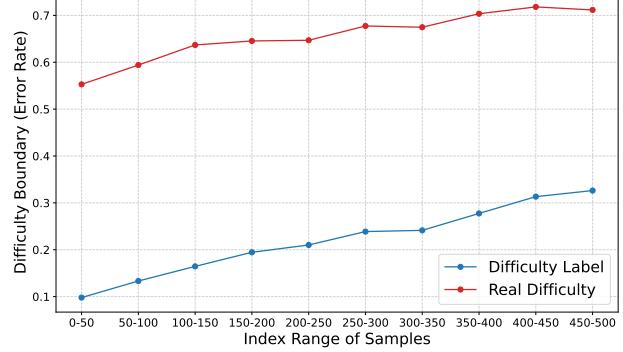


Figure 4. Trends of real and labeled difficulty over the index.

requires the model to engage in explicit reasoning, along with the adoption of prediction-label inconsistency as an objective metric.

Difficulty Boundary. With our proposed Difficulty Diffusion Mechanism and Difficulty Strategy Guidance, the difficulty boundary is significantly extended, as evidenced by an increase in error rate from 0.515 to 0.768, validating their effectiveness. Additionally, we analyze how the actual difficulty and difficulty labels evolve with the order of generated samples. As illustrated in Figure 4, both the difficulty label and actual difficulty exhibit a continuous

upward trend. This not only confirms that Difficulty Diffusion Mechanism operates as intended but also visually demonstrates the strong consistency between difficulty label and actual difficulty.

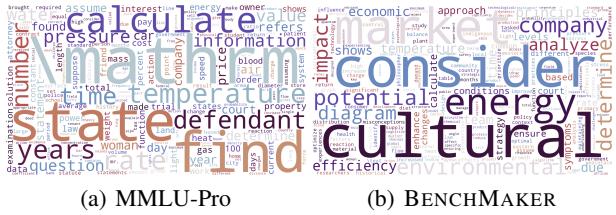


Figure 5. Word cloud of MMLU-Pro and the benchmark generated by BENCHMAKER under similar assessment demands.

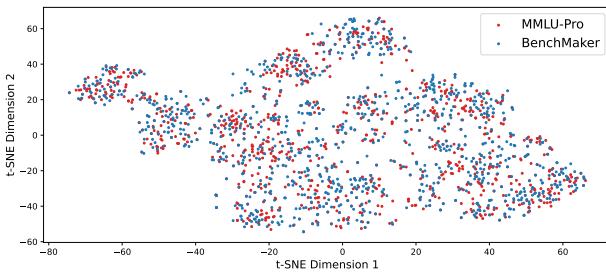


Figure 6. T-SNE results on the text embeddings of benchmarks.

5.3. A Closer Look at the Generated Benchmark

After metric analysis, we perform a more thorough examination of BENCHMAKER. Some of the generated samples are shown in Appendix H.

Lexical and Semantic. First, despite the obvious differences in word distribution between the generated benchmark and MMLU-Pro (Figure 5), it remains closely aligned with the domains covered by MMLU-Pro, demonstrating strong task alignment. Meanwhile, The semantic alignment between the two is more pronounced (Figure 6). Notably, the input demands (Appendix K) do not mention any information related to MMLU-Pro, effectively preventing the model from achieving a high degree of alignment by memorizing and replicating samples from MMLU-Pro.

Actual Error Rate. Although LLM-as-a-judge has provided an unbiased estimation of the benchmark’s faithfulness, we additionally conduct a manual check on 80 randomly selected samples. Our findings indicate that 3 samples have incorrect labels, 3 samples lack a correct candidate, resulting in an overall error rate of 7.5%. Meanwhile, LLM-as-a-judge identifies 5 problematic samples, with 3 overlapping with human judgment. These results suggest that: (1) BENCHMAKER still has room for improvement in faithfulness; (2) LLM-as-a-judge can serve as a partial proxy for human evaluation.

5.4. Reliability Estimation.

Since faithfulness of the generated benchmark cannot be totally ensured, we are curious about the effects of incorrect samples: Let \bar{a} and \bar{b} be the observed accuracies of two models A and B on the generated benchmark of size N , where a fraction K of the samples are incorrect (which can be estimated by the LLM-as-a-judge). Suppose that $\bar{a} > \bar{b}$, we aim to estimate the probability that the observed ability rank is correct (the true accuracies satisfy $E[a] > E[b]$). See detailed derivation in Appendix A. Suppose A and B have the same accuracy p on incorrect samples, we get:

$$E[a] = \frac{\bar{a} - K \cdot p}{1 - K}, \quad E[b] = \frac{\bar{b} - K \cdot p}{1 - K} \quad (4)$$

and

$$E[a] - E[b] = \frac{\bar{a} - \bar{b}}{1 - K} \quad (5)$$

Next, we perform hypothesis testing to assess the probability of $E[a] > E[b]$. We assume that $\bar{a} - \bar{b}$ follows a normal distribution. The z -score for the difference is:

$$\begin{aligned} z &= \frac{(\bar{a} - \bar{b}) / (1 - K)}{\sqrt{(\bar{a}(1 - \bar{a}) + \bar{b}(1 - \bar{b})) / (N(1 - K)^2)}} \\ &= \frac{(\bar{a} - \bar{b})\sqrt{N}}{\sqrt{\bar{a}(1 - \bar{a}) + \bar{b}(1 - \bar{b})}} \end{aligned} \quad (6)$$

where $\Phi(z)$ is the cumulative distribution function of the standard normal distribution. The probability $P(E[a] > E[b])$ is given by the right-tail probability of the normal distribution:

$$P(E[a] > E[b]) = 1 - \Phi(z) \quad (7)$$

where $\Phi(z)$ is the cumulative distribution function of the standard normal distribution. We can assess the reliability of BENCHMAKER evaluation results using (7). Also, we notice that K has the same scaling effect on both the numerator and denominator of the test statistic, thus does not alter the z -score. Consequently, as long as there is no bias, a certain proportion of noise in the benchmark will not affect the statistical significance of ability ranking.

Conclusions

The rapid advancement of large language models has driven an urgent demand for a generic benchmark generator. To this end, we first propose a comprehensive, automated, and unbiased evaluation framework to validate and optimize the reliability of benchmark generators. Based on this, we develop the BENCHMAKER method for reliable, generic, and efficient benchmark generation. Comprehensive experiments across multiple tasks and LLMs demonstrate that BENCHMAKER achieves human-aligned benchmark quality, with superior efficiency and generalization.

440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495

Impact Statement

This paper presents BENCHMAKER, an LLM-driven reliable, generic and efficient benchmark generator. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

Ethics Statement. All of the datasets used in this study were publicly available. Multiple authors jointly conducted the manual check for the Actual Error Rate section, and no extra annotators were employed for our data collection. We confirm that the datasets we used did not contain any harmful content and was consistent with their intended use (research). We have cited the datasets and relevant works used in this study.

References

- Anthropic. Claude 3.5. <https://www.anthropic.com/news/3-5-models-and-computer-use>.
- Balloccu, S., Schmidtová, P., Lango, M., and Dusek, O. Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source llms. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian's, Malta, March 17-22, 2024*, pp. 67–93. Association for Computational Linguistics, 2024. URL <https://aclanthology.org/2024.eacl-long.5>.
- Chan, X., Wang, X., Yu, D., Mi, H., and Yu, D. Scaling synthetic data creation with 1,000,000,000 personas. *CoRR*, abs/2406.20094, 2024. doi: 10.48550/ARXIV.2406.20094. URL <https://doi.org/10.48550/arXiv.2406.20094>.
- Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang, C., Wang, Y., Ye, W., Zhang, Y., Chang, Y., Yu, P. S., Yang, Q., and Xie, X. A survey on evaluation of large language models. *ACM Trans. Intell. Syst. Technol.*, 15(3):39:1–39:45, 2024. doi: 10.1145/3641289. URL <https://doi.org/10.1145/3641289>.
- Glazer, E., Erdil, E., Besiroglu, T., Chicharro, D., Chen, E., Gunning, A., Olsson, C. F., Denain, J., Ho, A., de Oliveira Santos, E., Järvinemi, O., Barnett, M., Sandler, R., Vrzala, M., Sevilla, J., Ren, Q., Pratt, E., Levine, L., Barkley, G., Stewart, N., Grechuk, B., Grechuk, T., Enugandla, S. V., and Wildon, M. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in AI. *CoRR*, abs/2411.04872, 2024. doi: 10.48550/ARXIV.2411.04872. URL <https://doi.org/10.48550/arXiv.2411.04872>.

Guo, X. and Vosoughi, S. Length does matter: Summary length can bias summarization metrics. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 15869–15879. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.984. URL <https://doi.org/10.18653/v1/2023.emnlp-main.984>.

Hamming, R. W. Error detecting and error correcting codes. *The Bell system technical journal*, 29(2):147–160, 1950.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021a. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.

Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the MATH dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, 2021b. URL <https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html>.

Huang, J., Chen, X., Mishra, S., Zheng, H. S., Yu, A. W., Song, X., and Zhou, D. Large language models cannot self-correct reasoning yet. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=Ikmd3fKBPQ>.

Hurst, A., Lerer, A., Goucher, A. P., Perelman, A., Ramesh, A., Clark, A., Ostrow, A., Welihinda, A., Hayes, A., Radford, A., et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

Ji, Z., Yu, T., Xu, Y., Lee, N., Ishii, E., and Fung, P. Towards mitigating LLM hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pp. 1827–1843. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-EMNLP.123. URL <https://doi.org/10.18653/v1/2023.findings-emnlp.123>.

Kaddour, J., Lynch, A., Liu, Q., Kusner, M. J., and Silva, R. Causal machine learning: A survey and open problems. *CoRR*, abs/2206.15475, 2022. doi: 10.48550/ARXIV.

- 495 2206.15475. URL <https://doi.org/10.48550/arXiv.2206.15475>.
- 496
- 497
- 498 Lei, F., Liu, Q., Huang, Y., He, S., Zhao, J., and Liu, K.
- 499 S3eval: A synthetic, scalable, systematic evaluation suite
- 500 for large language models. *CoRR*, abs/2310.15147, 2023.
- 501 doi: 10.48550/ARXIV.2310.15147. URL <https://doi.org/10.48550/arXiv.2310.15147>.
- 502
- 503 Lewis, P. S. H., Perez, E., Piktus, A., Petroni, F.,
- 504 Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih,
- 505 W., Rocktäschel, T., Riedel, S., and Kiela, D. Retrieval-
- 506 augmented generation for knowledge-intensive NLP
- 507 tasks. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information*
- 508 *Processing Systems 2020, NeurIPS 2020, December 6-12,*
- 509 *2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>.
- 510
- 511
- 512
- 513
- 514
- 515 Li, J., Hu, R., Huang, K., Zhuang, Y., Liu, Q., Zhu, M.,
- 516 Shi, X., and Lin, W. Perteval: Unveiling real knowledge
- 517 capacity of llms with knowledge-invariant perturbations.
- 518 *CoRR*, abs/2405.19740, 2024a. doi: 10.48550/ARXIV.2405.19740. URL <https://doi.org/10.48550/arXiv.2405.19740>.
- 519
- 520
- 521 Li, S., Huang, J., Zhuang, J., Shi, Y., Cai, X., Xu, M.,
- 522 Wang, X., Zhang, L., Ke, G., and Cai, H. Scilitllm:
- 523 How to adapt llms for scientific literature understanding.
- 524 *CoRR*, abs/2408.15545, 2024b. doi: 10.48550/ARXIV.2408.15545. URL <https://doi.org/10.48550/arXiv.2408.15545>.
- 525
- 526
- 527
- 528 Li, Z., Zhou, Z., Yao, Y., Li, Y., Cao, C., Yang, F., Zhang,
- 529 X., and Ma, X. Neuro-symbolic data generation for math
- 530 reasoning. *CoRR*, abs/2412.04857, 2024c. doi: 10.48550/ARXIV.2412.04857. URL <https://doi.org/10.48550/arXiv.2412.04857>.
- 531
- 532
- 533 Liu, Y., Iter, D., Xu, Y., Wang, S., Xu, R., and Zhu, C.
- 534 G-eval: NLG evaluation using gpt-4 with better human
- 535 alignment. In *Proceedings of the 2023 Conference*
- 536 *on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp.
- 537 2511–2522. Association for Computational Linguistics,
- 538 2023. URL <https://aclanthology.org/2023.emnlp-main.153>.
- 539
- 540
- 541
- 542 Long, L., Wang, R., Xiao, R., Zhao, J., Ding, X., Chen,
- 543 G., and Wang, H. On llms-driven synthetic data genera-
- 544 tion, curation, and evaluation: A survey. In *Findings*
- 545 *of the Association for Computational Linguistics, ACL*
- 546 *2024, Bangkok, Thailand and virtual meeting, August*
- 547 *11-16, 2024*, pp. 11065–11082. Association for Com-
- 548 putational Linguistics, 2024. doi: 10.18653/V1/2024.
- 549
- 550 FINDINGS-ACL.658. URL <https://doi.org/10.18653/v1/2024.findings-acl.658>.
- 551
- 552 Luo, H., Sun, Q., Xu, C., Zhao, P., Lin, Q., Lou, J.,
- 553 Chen, S., Tang, Y., and Chen, W. Arena learning:
- 554 Build data flywheel for llms post-training via simulated
- 555 chatbot arena. *CoRR*, abs/2407.10627, 2024a. doi:
- 556 10.48550/ARXIV.2407.10627. URL <https://doi.org/10.48550/arXiv.2407.10627>.
- 557
- 558 Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W.,
- 559 Tao, C., Ma, J., Lin, Q., and Jiang, D. Wizardcoder:
- 560 Empowering code large language models with evol-
- 561 instruct. In *The Twelfth International Conference on*
- 562 *Learning Representations, ICLR 2024, Vienna, Austria,*
- 563 *May 7-11, 2024*. OpenReview.net, 2024b. URL <https://openreview.net/forum?id=UnUwSIgK5W>.
- 564
- 565 Maheshwari, G., Ivanov, D., and Haddad, K. E. Efficacy of
- 566 synthetic data as a benchmark. *CoRR*, abs/2409.11968,
- 567 2024. doi: 10.48550/ARXIV.2409.11968. URL <https://doi.org/10.48550/arXiv.2409.11968>.
- 568
- 569 Montahaei, E., Alihosseini, D., and Baghshah, M. S. Jointly
- 570 measuring diversity and quality in text generation models.
- 571 *CoRR*, abs/1904.03971, 2019. URL <http://arxiv.org/abs/1904.03971>.
- 572
- 573 OpenAI. text-embedding-ada-002. <https://platform.openai.com/docs/guides/embeddings>.
- 574
- 575 Perez, E., Ringer, S., Lukosiuete, K., Nguyen, K., Chen,
- 576 E., Heiner, S., Pettit, C., Olsson, C., Kundu, S., Kadavath, S., Jones, A., Chen, A., Mann, B., Israel, B., Seethor, B., McKinnon, C., Olah, C., Yan, D., Amodei, D., Amodei, D., Drain, D., Li, D., Tran-Johnson, E., Khundadze, G., Kernion, J., Landis, J., Kerr, J., Mueller, J., Hyun, J., Landau, J., Ndousse, K., Goldberg, L., Lovitt, L., Lucas, M., Sellitto, M., Zhang, M., Kingsland, N., Elhage, N., Joseph, N., Mercado, N., DasSarma, N., Rausch, O., Larson, R., McCandlish, S., Johnston, S., Kravec, S., Showk, S. E., Lanham, T., Telleen-Lawton, T., Brown, T., Henighan, T., Hume, T., Bai, Y., Hatfield-Dodds, Z., Clark, J., Bowman, S. R., Askell, A., Grosse, R., Hernandez, D., Ganguli, D., Hubinger, E., Schiefer, N., and Kaplan, J. Discovering language model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pp. 13387–13434. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-ACL.847. URL <https://doi.org/10.18653/v1/2023.findings-acl.847>.
- 577
- 578 Tam, Z. R., Wu, C., Tsai, Y., Lin, C., Lee, H., and Chen, Y.
- 579 Let me speak freely? A study on the impact of format

- 550 restrictions on large language model performance. In *Proceedings of the 2024 Conference on Empirical Methods in*
 551 *Natural Language Processing: EMNLP 2024 - Industry*
 552 *Track, Miami, Florida, USA, November 12-16, 2024*, pp.
 553 1218–1236. Association for Computational Linguistics,
 554 2024. URL <https://aclanthology.org/2024.emnlp-industry.91>.
- 555 Tao, Z., Lin, T., Chen, X., Li, H., Wu, Y., Li, Y.,
 556 Jin, Z., Huang, F., Tao, D., and Zhou, J. A
 557 survey on self-evolution of large language models.
 558 *CoRR*, abs/2404.14387, 2024. doi: 10.48550/ARXIV.
 559 2404.14387. URL <https://doi.org/10.48550/arXiv.2404.14387>.
- 560 Thakur, A. S., Choudhary, K., Ramayapally, V. S.,
 561 Vaidyanathan, S., and Hupkes, D. Judging the judges:
 562 Evaluating alignment and vulnerabilities in llms-as-judges.
 563 *CoRR*, abs/2406.12624, 2024. doi: 10.48550/ARXIV.2406.12624. URL <https://doi.org/10.48550/arXiv.2406.12624>.
- 564 Vallat, R. Pingouin: statistics in python. *J. Open*
 565 *Source Softw.*, 3(31):1026, 2018. doi: 10.21105/
 566 JOSS.01026. URL <https://doi.org/10.21105/joss.01026>.
- 567 Wang, K., Zhu, J., Ren, M., Liu, Z., Li, S., Zhang, Z.,
 568 Zhang, C., Wu, X., Zhan, Q., Liu, Q., and Wang, Y.
 569 A survey on data synthesis and augmentation for large
 570 language models. *CoRR*, abs/2410.12896, 2024a. doi:
 571 10.48550/ARXIV.2410.12896. URL <https://doi.org/10.48550/arXiv.2410.12896>.
- 572 Wang, X., Wei, J., Schuurmans, D., Le, Q. V., Chi,
 573 E. H., Narang, S., Chowdhery, A., and Zhou, D.
 574 Self-consistency improves chain of thought reasoning
 575 in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net,
 576 2023a. URL <https://openreview.net/forum?id=1PL1NIMMrw>.
- 577 Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A.,
 578 Khashabi, D., and Hajishirzi, H. Self-instruct: Aligning
 579 language models with self-generated instructions. In
 580 *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pp.
 581 13484–13508. Association for Computational Linguistics, 2023b. doi: 10.18653/V1/2023.ACL-LONG.754. URL <https://doi.org/10.18653/v1/2023.acl-long.754>.
- 582 Wang, Y., Ma, X., Zhang, G., Ni, Y., Chandra, A., Guo, S.,
 583 Ren, W., Arulraj, A., He, X., Jiang, Z., Li, T., Ku, M.,
 584 Wang, K., Zhuang, A., Fan, R., Yue, X., and Chen, W.
 585 Mmlu-pro: A more robust and challenging multi-task lan-
 586 guage understanding benchmark. *CoRR*, abs/2406.01574,
 587 2024b. doi: 10.48550/ARXIV.2406.01574. URL <https://doi.org/10.48550/arXiv.2406.01574>.
- 588 Wu, S., Huang, Y., Gao, C., Chen, D., Zhang, Q., Wan, Y.,
 589 Zhou, T., Zhang, X., Gao, J., Xiao, C., and Sun, L. Unigen:
 590 A unified framework for textual dataset generation
 591 using large language models. *CoRR*, abs/2406.18966,
 592 2024. doi: 10.48550/ARXIV.2406.18966. URL <https://doi.org/10.48550/arXiv.2406.18966>.
- 593 Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C.,
 594 Li, C., Li, C., Liu, D., Huang, F., Dong, G., Wei, H., Lin,
 595 H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J.,
 596 Yang, J., Xu, J., Zhou, J., Bai, J., He, J., Lin, J., Dang,
 597 K., Lu, K., Chen, K., Yang, K., Li, M., Xue, M., Ni, N.,
 598 Zhang, P., Wang, P., Peng, R., Men, R., Gao, R., Lin,
 599 R., Wang, S., Bai, S., Tan, S., Zhu, T., Li, T., Liu, T.,
 600 Ge, W., Deng, X., Zhou, X., Ren, X., Zhang, X., Wei,
 601 X., Ren, X., Liu, X., Fan, Y., Yao, Y., Zhang, Y., Wan,
 602 Y., Chu, Y., Liu, Y., Cui, Z., Zhang, Z., Guo, Z., and
 603 Fan, Z. Qwen2 technical report. *CoRR*, abs/2407.10671,
 604 2024. doi: 10.48550/ARXIV.2407.10671. URL <https://doi.org/10.48550/arXiv.2407.10671>.
- 605 Yu, L., Jiang, W., Shi, H., Yu, J., Liu, Z., Zhang, Y., Kwok,
 606 J. T., Li, Z., Weller, A., and Liu, W. Metamath: Boot-
 607 strap your own mathematical questions for large language
 608 models. In *The Twelfth International Conference on
 609 Learning Representations, ICLR 2024, Vienna, Austria,
 610 May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=N8N0hgNDrt>.
- 611 Yu, Y., Zhuang, Y., Zhang, J., Meng, Y., Ratner, A. J.,
 612 Krishna, R., Shen, J., and Zhang, C. Large language
 613 model as attributed training data generator: A tale
 614 of diversity and bias. In *Advances in Neural Infor-
 615 mation Processing Systems 36: Annual Conference on
 616 Neural Information Processing Systems 2023, NeurIPS
 617 2023, New Orleans, LA, USA, December 10 - 16, 2023*. URL http://papers.nips.cc/paper_files/paper/2023/hash/ae9500c4f5607caf2eff033c67daa9d7-Abstract-Dataset_and_Benchmarks.html.
- 618 Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi,
 619 Y. Hellaswag: Can a machine really finish your sentence?
 620 In *Proceedings of the 57th Conference of the Association
 621 for Computational Linguistics, ACL 2019, Florence, Italy,
 622 July 28- August 2, 2019, Volume 1: Long Papers*, pp.
 623 4791–4800. Association for Computational Linguistics,
 624 2019. doi: 10.18653/V1/P19-1472. URL <https://doi.org/10.18653/v1/p19-1472>.
- 625 Zheng, L., Chiang, W., Sheng, Y., Zhuang, S., Wu,
 626 Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. P.,

- 605 Zhang, H., Gonzalez, J. E., and Stoica, I. Judging
606 llm-as-a-judge with mt-bench and chatbot arena. In
607 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets_and_Benchmarks.html.
- 614
- 615 Zhu, K., Chen, J., Wang, J., Gong, N. Z., Yang, D., and
616 Xie, X. Dyval: Dynamic evaluation of large language
617 models for reasoning tasks. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net,
618 2024a. URL <https://openreview.net/forum?id=gjfOL9z5Xr>.
- 622
- 623 Zhu, K., Wang, J., Zhao, Q., Xu, R., and Xie, X. Dyval 2:
624 Dynamic evaluation of large language models by meta
625 probing agents. *CoRR*, abs/2402.14865, 2024b. doi:
626 10.48550/ARXIV.2402.14865. URL <https://doi.org/10.48550/arXiv.2402.14865>.
- 628
- 629 Zhu, Y., Lu, S., Zheng, L., Guo, J., Zhang, W., Wang,
630 J., and Yu, Y. Texxygen: A benchmarking platform
631 for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pp. 1097–1100. ACM, 2018.
632 doi: 10.1145/3209978.3210080. URL <https://doi.org/10.1145/3209978.3210080>.
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647
- 648
- 649
- 650
- 651
- 652
- 653
- 654
- 655
- 656
- 657
- 658
- 659

660 A. Derivation of the Probability for $E[a] > E[b]$

661 In this appendix, we present a detailed derivation of the probability that the expected accuracy $E[a]$ of model A is greater
 662 than the expected accuracy $E[b]$ of model B in a noiseless benchmark, given the noisy benchmark observations. We will use
 663 the following notation throughout:

- 664 • N : The total number of samples in the benchmark.
- 665 • K : The proportion of samples with incorrect labels, i.e., the noise ratio, where $K \in [0, 1]$.
- 666 • \bar{a} : The observed accuracy of model A on the noisy benchmark.
- 667 • \bar{b} : The observed accuracy of model B on the noisy benchmark.
- 668 • p : The probability that both models predict the incorrect label correctly on a noisy sample. This probability is assumed
 669 to be identical for both models on incorrect labels.
- 670 • $E[a]$: The expected accuracy of model A on the noiseless benchmark.
- 671 • $E[b]$: The expected accuracy of model B on the noiseless benchmark.

672 Given these notations, our goal is to determine the probability that $E[a] > E[b]$ based on the noisy observed accuracies \bar{a}
 673 and \bar{b} .

674 Step 1: Relating \bar{a} and \bar{b} to $E[a]$ and $E[b]$

675 From the given setup, we know that the observed accuracies \bar{a} and \bar{b} can be written as a weighted average of the expected
 676 accuracies $E[a]$ and $E[b]$ on the correct samples, and the probability p on the incorrect samples. Specifically, the formulas
 677 for \bar{a} and \bar{b} are:

$$\begin{aligned} 678 \bar{a} &= (1 - K) \cdot E[a] + K \cdot p \\ 679 \bar{b} &= (1 - K) \cdot E[b] + K \cdot p \end{aligned}$$

680 These equations express the observed accuracy of each model as the weighted average of the correct label samples and
 681 the noisy (incorrect label) samples. The weight $(1 - K)$ represents the proportion of correct labels, and K represents the
 682 proportion of incorrect labels.

683 Step 2: Solving for $E[a]$ and $E[b]$

684 To isolate $E[a]$ and $E[b]$, we rearrange the above equations:

$$\begin{aligned} 685 E[a] &= \frac{\bar{a} - K \cdot p}{1 - K} \\ 686 E[b] &= \frac{\bar{b} - K \cdot p}{1 - K} \end{aligned}$$

687 Thus, $E[a]$ and $E[b]$ are directly related to the observed accuracies \bar{a} and \bar{b} , and the noise ratio K .

688 Step 3: Comparing $E[a]$ and $E[b]$

689 To determine the probability that $E[a] > E[b]$, we first compute the difference between $E[a]$ and $E[b]$:

$$690 E[a] - E[b] = \frac{\bar{a} - K \cdot p}{1 - K} - \frac{\bar{b} - K \cdot p}{1 - K}$$

691 Simplifying this expression:

$$E[a] - E[b] = \frac{\bar{a} - \bar{b}}{1 - K}$$

Thus, $E[a] > E[b]$ if and only if $\bar{a} - \bar{b} > 0$, which indicates that the observed accuracy of model A must be greater than that of model B in the noisy benchmark for $E[a]$ to exceed $E[b]$ in the noiseless benchmark.

Step 4: Statistical Hypothesis Testing

To quantify the probability of $E[a] > E[b]$, we perform a hypothesis test on $\bar{a} - \bar{b}$, assuming that both \bar{a} and \bar{b} are derived from binomial distributions (since they represent the correct classification probabilities on the noisy benchmark). We assume that the difference $\bar{a} - \bar{b}$ follows a normal distribution under certain conditions (via the Central Limit Theorem). Thus, the expected value of $\bar{a} - \bar{b}$ is:

$$\mathbb{E}[\bar{a} - \bar{b}] = \frac{\bar{a} - \bar{b}}{1 - K}$$

The variance of $\bar{a} - \bar{b}$, assuming independent samples, is given by:

$$\begin{aligned} \text{Var}[\bar{a} - \bar{b}] &= \frac{\text{Var}[\bar{a}] + \text{Var}[\bar{b}]}{(1 - K)^2} \\ &= \frac{\bar{a}(1 - \bar{a})/N + \bar{b}(1 - \bar{b})/N}{(1 - K)^2} \end{aligned} \tag{8}$$

where $\text{Var}[\bar{a}]$ and $\text{Var}[\bar{b}]$ are the variances of the observed accuracies of models A and B, respectively. These variances can be computed from the binomial distributions underlying \bar{a} and \bar{b} .

Now, we define the z -score for the observed difference $\bar{a} - \bar{b}$ as follows:

$$\begin{aligned} z &= \frac{\mathbb{E}[\bar{a} - \bar{b}]}{\sqrt{\text{Var}[\bar{a} - \bar{b}]}} \\ &= \frac{(\bar{a} - \bar{b})/(1 - K)}{\sqrt{(\bar{a}(1 - \bar{a}) + \bar{b}(1 - \bar{b}))/N(1 - K)^2}} \\ &= \frac{(\bar{a} - \bar{b})\sqrt{N}}{\sqrt{\bar{a}(1 - \bar{a}) + \bar{b}(1 - \bar{b})}} \end{aligned} \tag{9}$$

This z -score follows a standard normal distribution. The probability that $E[a] > E[b]$ is the probability that $\bar{a} - \bar{b} > 0$, which is equivalent to:

$$P(\bar{a} - \bar{b} > 0) = P(z > 0) = 1 - \Phi(z)$$

where $\Phi(z)$ is the cumulative distribution function (CDF) of the standard normal distribution. Thus, the p -value is given by:

$$p\text{-value} = 1 - \Phi(z)$$

Step 5: Conclusion

To summarize, the probability that $E[a] > E[b]$ in the noiseless benchmark, given the noisy benchmark observations, is determined by the observed accuracy difference $\bar{a} - \bar{b}$ and the noise ratio K . The probability is computed using a hypothesis test on $\bar{a} - \bar{b}$, assuming it follows a normal distribution. The final formula for this probability is:

$$P(E[a] > E[b]) = 1 - \Phi(z)$$

where $\Phi(z)$ is the CDF of the standard normal distribution and z is the computed z -score. This result allows us to assess the likelihood that model A has a higher expected accuracy than model B in the noiseless benchmark based on noisy observations.

B. Unsuccessful Attempts for Optimizing Benchmark Generator

B.1. Faithfulness

We explored the widely studied self-correction strategy to improve the faithfulness of benchmarks. Specifically, for each generated sample, the model first acts as a judge and then refines samples it deems insufficiently faithful. However, our preliminary results indicate that while this approach yields minor improvements in mathematical tasks, it provides little benefit for tasks such as MMLU-Pro and instead introduces additional computational overhead.

B.2. Difficulty Controllability

As previously mentioned, we attempted to have the model generate samples with specified difficulty levels, but the resulting samples exhibited low difficulty differentiation. To address this, we further explored having the model assess the difficulty of its generated samples. However, this strategy yielded promising results only on the MATH task.

B.3. Difficulty Diffusion Mechanism

Previous studies (Wang et al., 2024b) have attempted to increase question difficulty by expanding the number of answer choices. However, our experiments show that scaling up the number of candidates quickly reaches a saturation point. We hypothesize that this is due to the model’s difficulty in generating a large number of sufficiently deceptive distractors.

B.4. Diversity

To enhance sample diversity, in addition to AttrPrompt, we experimented with assigning different personas (Chan et al., 2024) to the model and instructing it to generate characteristic samples based on its assigned persona. However, we found that this approach was not particularly effective for the MATH task, especially in semantic diversity.

C. Data from Huggingface

We obtained information on open-source model releases and download counts from the Hugging Face API (`from huggingface_hub import HfApi`). Since the number of open-source model releases far exceeds that of closed-source models, we use the former to represent the “Number of Language Model Releases.” Additionally, as Hugging Face does not provide monthly download counts for each model, we use the historical total downloads of models released within a given statistical period as the total downloads for that period. The corresponding code is shown below.

D. Difficulty Levels

- **Level 1:** The simplest, equivalent to lower-grade elementary school
- **Level 2:** Relatively simple, equivalent to upper-grade elementary school
- **Level 3:** Simple, equivalent to middle school
- **Level 4:** Average, equivalent to high school
- **Level 5:** Slightly difficult, equivalent to university student
- **Level 6:** Difficult, equivalent to Master’s
- **Level 7:** Quite difficult, equivalent to PhD student

- **Level 8:** Very difficult, equivalent to professor
- **Level 9:** Extremely difficult, equivalent to field expert
- **Level 10:** Most difficult, equivalent to top human level or beyond human level

E. Benchmarking Model List

- **phoenix-inst-chat-7b:** <https://huggingface.co/FreedomIntelligence/phoenix-inst-chat-7b>
- **vicuna-7b-v1.3:** <https://huggingface.co/lmsys/vicuna-7b-v1.3>
- **Qwen2.5-3B:** <https://huggingface.co/Qwen/Qwen2.5-3B>
- **phi-2:** <https://huggingface.co/microsoft/phi-2>
- **Phi-3.5-mini-instruct:** <https://huggingface.co/microsoft/Phi-3.5-mini-instruct>
- **Yi-1.5-6B-Chat:** <https://huggingface.co/01-ai/Yi-1.5-6B-Chat>
- **Qwen2.5-7B:** <https://huggingface.co/Qwen/Qwen2.5-7B>
- **vicuna-7b-v1.5:** <https://huggingface.co/lmsys/vicuna-7b-v1.5>
- **Qwen2-1.5B-Instruct:** <https://huggingface.co/Qwen/Qwen2-1.5B-Instruct>
- **phoenix-inst-chat-7b-v1.1:** <https://huggingface.co/FreedomIntelligence/phoenix-inst-chat-7b-v1.1>
- **Qwen-Plus:** <https://huggingface.co/Qwen>
- **GPT-3.5 turbo:** <https://openai.com/index/gpt-3-5-turbo-fine-tuning-and-api-updates/>

F. Details of Difficulty Diffusion Mechanism

Given that the LLM has a certain level of difficulty perception, we iteratively select the more challenging samples according to β from the generated ones as difficulty references, and instruct the LLM to generate a more difficult sample. Specifically, To prevent reference samples from becoming overly fixed, which may lead to homogenization in generated samples, we adopt the following strategy:

1. We track the number of times each sample x_i has been used as a reference sample, denoted as t_i , and compute a calibrated difficulty label:

$$\text{Calibrate_Difficulty} = \text{Difficulty_Label} \times 0.9^{t_i/\text{Reference_Number}} \quad (10)$$

The samples are then sorted based on this adjusted difficulty.

2. Each time, we select $2 \times \text{Reference_Number}$ samples with the highest Calibrate_Difficulty as candidates. From this pool, we randomly sample Reference_Number as reference samples and shuffle their order.

Our preliminary experiments indicate a positive correlation between problem difficulty and Reference_Number. In our experiments, we set Reference_Number to 8. This allows the sample difficulty to rise continuously through diffusion.

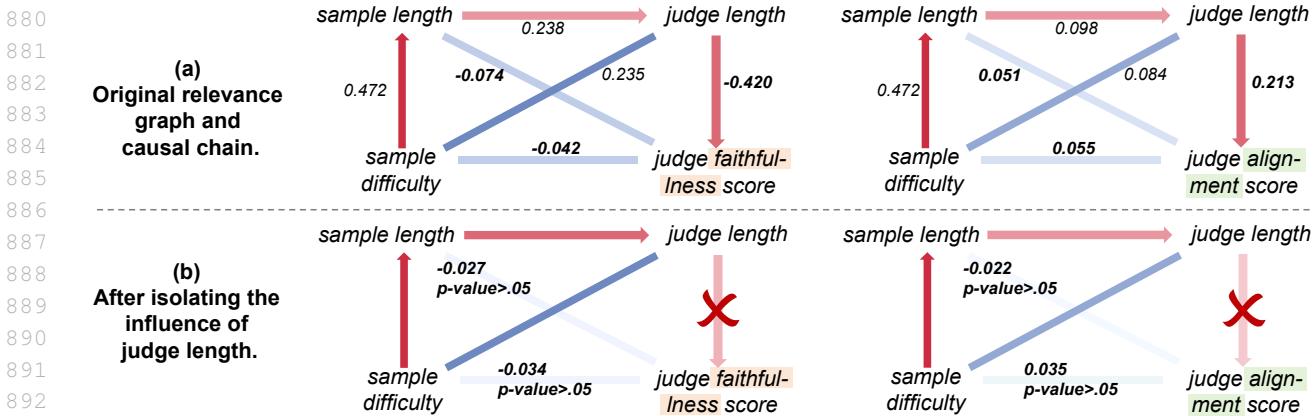


Figure 7. Pearson correlations among key factors of benchmark evaluation and LLM (GPT-4o mini) judge scores (faithfulness and alignment). The most relevant path of each subject is highlighted in red to show the possible causal chain.

G. Converting Benchmark Sample Format

MCQ to OTG Format. By removing the options from the samples and using the solution and answer corresponding to the correct option as the ground truth, we can easily transform the MCQ-style benchmark generated based on MATH assessment demands into an open-ended text generation (OTG) benchmark. Comparing these two benchmark formats, we find that the OTG format makes the questions more challenging (error rate: $0.865 > 0.768$) and results in lower knowledge diversity (hamming distance: $0.365 < 0.403$). We attribute this to the model’s inability to rely on option cues to answer certain questions, which leads to a large portion of the knowledge vector being zero, thereby reducing knowledge diversity. Additionally, we observe a decline in the benchmark’s effectiveness (pearson: $0.915 < 0.935$), which we hypothesize is indirectly caused by the drop in knowledge diversity.

OTG to MCQ Format. To convert the MATH benchmark into a MCQ format, we employed GPT-4o, GPT-4o mini, GPT-3.5 turbo, and Claude-3.5 Haiku to answer MATH problems, sampling 10 responses per model with a temperature of 1. We identified the three most frequently occurring incorrect answers as distractors while retaining the correct answer’s rationale. If the number of incorrect answers was insufficient, we supplemented them using GPT-4o mini. Through this process, we successfully transformed the MATH benchmark from an OTG format to an MCQ format.

H. Examples of the Generated Samples

MATH:

Example 1:
A researcher is studying the distribution of three specific proteins in a cell. There are 4 locations within the cell where each protein can be present. However, due to experimental conditions, at least one protein must be present in each location. In how many different ways can the proteins be distributed in the cell, considering overlap in presence is allowed?

- A. 2187
- B. 2401
- C. 4096
- D. 2048

Label:B

Example 2:

Find the smallest positive integer n such that n is divisible by 6, 10, and 15, and $n \equiv 2 \pmod{4}$.

- A. 120
- B. 20
- C. 90
- D. 30

935
936 Label:D
937

938 MMLU-Pro:

939 Example 1:

940 A 45-year-old woman with type 2 diabetes decides to improve her health by adopting a low-carbohydrate, high-protein diet, starting a daily 30-minute brisk walk routine, and taking a new medication that increases insulin sensitivity. She also begins consuming a herbal supplement believed to enhance energy levels. After two months, she notices an increase in fatigue, frequent headaches, and unexplained weight gain. What is the most likely reason for her symptoms?

- 941 A. Low-carbohydrate diet leading to nutritional deficiencies
- 942 B. Brisk walk routine causing excessive physical exertion
- 943 C. Medication side effects causing insulin fluctuations
- 944 D. Herbal supplement causing hormonal imbalance
- 945 E. Increased protein intake causing kidney strain
- 946 F. Inadequate hydration from dietary changes
- 947 G. Overconsumption of high-protein foods leading to weight gain
- 948 H. Lack of fiber intake affecting metabolism
- 949 I. Decrease in carbohydrate intake causing energy depletion
- 950 J. Stress from lifestyle changes impacting health

951 Label:C

952 Example 2:

953 An architect is designing a complex apartment building, which features a series of irregularly shaped balconies. The layout of one of the building's wings is depicted in the accompanying diagram. Each balcony's area is defined by the function $f(x) = 3x^2 + 2x + 1$ over the interval $[0, 2]$ meters, representing a horizontal cross-section. The total length of the wing is 10 meters, and each balcony occurs at every meter along this length, aligned perpendicularly. To meet safety regulations, the architect needs to ensure that the probability of a randomly selected balcony having an area greater than 8 square meters is at least 0.5. Calculate the probability that a randomly selected balcony from this wing has an area greater than 8 square meters, using integration to determine the areas and probabilities involved. Consider potential pitfalls like incorrect integral setup or probability interpretation.

- 954 A. 0.1
- 955 B. 0.2
- 956 C. 0.3
- 957 D. 0.4
- 958 E. 0.5
- 959 F. 0.6
- 960 G. 0.7
- 961 H. 0.8
- 962 I. 0.9
- 963 J. 1.0

964 Label:J

965 HellaSwag:

966 During a family reunion, Mark is honored with the 'Outstanding Contributor' award for his recent volunteer efforts in the community. As he stands in front of his relatives, he expresses heartfelt gratitude towards everyone who supported him but fails to mention his younger sister, Lily, who organized the charity event that helped him earn this recognition. After the ceremony, Lily watches Mark celebrate with others, her face a mix of pride and disappointment. When Mark approaches her, excitedly asking, 'Did you see me win? I couldn't have done it without your help!' Given Lily's conflicting feelings about being overlooked, how is she most likely to respond?

- 967 A. 'Congratulations, Mark! I'm really proud of you! But I can't help feeling a bit overshadowed since I organized the event.'

```
990 B. 'Wow, Mark! You totally deserve this! Yet, it's tough for me to celebrate when my  
991 efforts went unnoticed.'  
992 C. 'That was an amazing award, Mark! I'm happy for you! However, it stings that my  
993 contribution was overlooked.'  
994 D. 'I'm so thrilled for you, Mark! Your achievement is incredible! But it feels a little  
995 unfair that I didn't get a shoutout.'
```

```
996 Label:C
```

```
997
```

```
998
```

I. Prompt List

```
1000
```

LLM as Faithfulness Judge:

```
1001 You are an expert who excels at analyzing whether a given response correctly answers a  
100 provided question.
```

```
100
```

```
100 **Question:**  
100 {{question}}
```

```
100
```

```
100 **Response to be Checked:**  
100 {{response}}
```

```
100
```

```
101 Please note that the given question may be unsolvable, have a unique solution, multiple  
101 solutions, etc.
```

```
101
```

```
Therefore, you should carefully analyze the correctness of the response to be checked  
101 based on the given question.
```

```
101
```

```
Here are the rules to strictly follow when analyzing the correctness of a response:
```

```
101
```

1. **Step-by-Step Analysis**: Analyze the response step by step, reviewing the reasoning and correctness of each step. For every step, first **restate and summarize** the reasoning logic and conclusion presented in the response, then analyze the correctness of that specific step.
2. **Focus on Evaluation**: Remember that your primary mission is to determine whether the reasoning process is correct. Avoid attempting to solve the problem yourself. Instead, focus strictly on analyzing the correctness of the response's reasoning process, one step at a time.
3. **Avoid Premature Judgments**: Do not rush to make judgments (such as claiming the response is flawed or completely correct) at the beginning. Ensure your evaluation is based on a thorough step-by-step analysis before arriving at a conclusion.
4. **Reverse Validation**: After completing the step-by-step analysis, substitute the answer back into the original problem and perform reverse validation of the parameters to cross-verify the correctness of the response.

```
After completing your analysis, please provide your judgment on the correctness of the  
102 response, as well as your confidence level in that judgment.
```

```
102
```

```
Your output should follow the template and example below:
```

```
103
```

```
Analyses:{Your detailed analyses}  
Judgement:{0: You think both the final answer of the response is wrong; 0.5: You think the  
103 reasoning path has some mistakes, but the final answer of the response is correct; 1:  
103 You agree with the reasoning path and the final answer of the response}
```

```
103
```

```
##Example##
```

```
103
```

```
Analyses:{Your detailed analyses}
```

```
103
```

```
Judgement:1
```

```
103
```

```
##Example End##
```

```
103
```

```
Now begin with "Analyses:"
```

```
1039
```

LLM as Comparison-based Faithfulness Judge:

```
1041
```

```
You are a knowledgeable expert with the task of analyzing the quality of a given question  
104 and its candidate answers.
```

```
104
```

```

104     ###Question
104     {{question}}
104
104     ###Candidate 1:
104     {{can1}}
104
105     ###Candidate 2:
105     {{can2}}
105
105     ###Your task: Correctness Analysis
105     1. Analyze whether the question is correct, reasonable, and clearly stated.
105     2. For the given question, analyze whether the provided ###Candidate 1 and ###Candidate 2
105        are correct step by step sequentially.
105        (Do not favor a candidate just because it is long; evaluate candidates strictly based on
105          correctness.)
105     3. Based on the above analysis, output your judgment of the question quality according to
105        the following scale:
105        0 point indicate an incorrect question with ambiguities and no uniquely suitable
105          answer among the options.
106        0.5 point indicates a minor error in the question, but there is still a uniquely
106          suitable answer among the options.
106        1 point indicate no errors in the question, with one uniquely correct answer among the
106          options.
106     4. Please also output your chosen correct option
106 You should follow the template below to output:
106 "##Faithfulness:{{score}}##, ##Label:{{}}##" (e.g., ##Faithfulness:2##, ##Label:B##).
106 Please note that if you believe there is no correct option or there are multiple correct
106 options, output ##Faithfulness:0##, ##Label:None##.
106
106 You should begin your response with "Correctness Analysis".
107

```

LLM as Relevance Judge:

```

107 You are an expert who excels at analyzing whether a given question can be used to assess a
107 specific ability.
107
107 **Question:***
107 {{question}}
107
107 **Ability:***
107 {{ability}}
107
108 You should first carefully analyze what abilities the given question can be used to test.
108 Based on this analysis, compare it with the given abilities.
108 After completing your analysis, please provide your judgment on whether the given question
108 can be used to test the given ability, as well as your confidence in that judgment.
108
108 Your output should follow the template below:
108 Analyses:{Your detailed analyses}
108 Judgement:
108 {output 0 if: You believe the given question is completely unable to test the given
108   ability;
108 output 0.5 if: You believe the given question is primarily meant to test other abilities,
108   but can also test the given ability to some extent;
108 output 1 if: You believe the given question primarily tests the given ability.}
108
109 Now begin with "Analyses:"
109

```

1094
 1095 **Directly Prompting LLM as Generic Benchmark Generator:** Notably, before allowing the LLM to formally generate
 1096 the benchmark, we first require it to produce descriptions for each part of the sample based on the assessment demands,
 1097 including Task Description, Query Description, and Option Description. This helps the model better understand and align
 1098 with the assessment demands, ensuring higher-quality and more consistent benchmark generation.
 1099

```

1100 You are a knowledgeable benchmark creator.
1101 Your task is to generate a creative questions based on the provided Task Description,
1102 Query Description, Option Description, Generation Guidelines, and Output Description
1103 to help build a benchmark that assesses the given task.
1104
1105 #### Task Description:
1106 {{task define}}
1107
1108 #### Query Description:
1109 {{query define}}
1110
1111 #### Option Description:
1112 {{option define}}
1113
1114 #### Generation Guidelines:
1115 1. Analyze the given task and think step-by-step about the content needed to construct the
1116     question, begin with "Analyses:".
1117 2. Generate the question content, begin with "Question:".
1118 3. Generate 10 candidates, with only one as the right option. Begin with "Candidates:".
1119 4. Generate the index of the right option, begin with "Right Option:".
1120
1121 #### Output Description:
1122 Strictly follow the template below to generate your sample.
1123 **Template**
1124 ##Analyses## {{You analyze the provided attributes and outline the process for
1125     constructing the question to be generated.}}
1126 ##Question## {{Your generated question content}}
1127 ##Candidates##
1128 {{Your generated Candidates}}
1129 ##Right Option##{{Index of the right option, e.g., B}}
1130 **Template End**
1131
1132 Attention: You need to **strictly follow the template** and don't generate any other
1133     contents. Begin your response with "##Analyses## "
1134
1135

```

J. Examples of the Generated Difficulty Strategies

MATH:

```

1136 Strategy 1:
1137 Complexity of Biological Concept is Basic
1138 Complexity of Biological Concept is Intermediate
1139 Complexity of Biological Concept is Advanced
1140
1141 Strategy 2:
1142 Required Reasoning Steps set as Single-step
1143 Required Reasoning Steps set as Multi-step (2-3 steps)
1144 Required Reasoning Steps set as Multi-step (4-6 steps)
1145 Required Reasoning Steps set as More than 6 steps
1146
1147 Strategy 3:
1148 Familiarity with the Topic is Common
1149 Familiarity with the Topic is Uncommon
1150 Familiarity with the Topic is Rare
1151
1152 Strategy 4:
1153 Type of Biological Data Analysis is Qualitative
1154 Type of Biological Data Analysis is Quantitative
1155 Type of Biological Data Analysis is Advanced Data Interpretation
1156
1157 Strategy 5:
1158 Application of Concepts is Direct
1159 Application of Concepts is Modified
1160

```

```

115 Application of Concepts is Novel
115
115 Strategy 6:
115 Integration Across Biological Disciplines is Single-discipline
115 Integration Across Biological Disciplines is Cross-disciplinary
115 Integration Across Biological Disciplines is Interdisciplinary
116
116 Strategy 7:
116 Depth of Required Knowledge is Surface-level
116 Depth of Required Knowledge is In-depth
116 Depth of Required Knowledge is Comprehensive
116

```

Prompt of BENCHMAKER:

```

116 You are a knowledgeable benchmark creator.
116 Your task is to generate a creative question based on the provided Task Description, Query
116     Description, Option Description, General Attributes Descriptions, Difficulty
116     Strategies Description, Generation Guidelines, and Output Description to help build a
116     benchmark that assesses the given task.
117
117     ### Overall Task Description:
117     {{original task}}
118
117     ### Detailed Task Description:
117     {{task define}}
118
117     ### Query Description:
117     {{query define}}
118
118     ### Option Description:
118     {{option define}}
118
118     ### General Attributes Description:
118 You can refer to the following attributes and their corresponding values to construct
118     questions, which means the questions you generate should ideally align with some of
118     these attributes.
118 Please note, if you find any conflicting or confusing parts among the attributes listed,
118     you may disregard them.
119     {{attribute define}}
119
119     ### Difficulty Strategies Description:
119 Your generated questions should meet the following difficulty attribute requirements. If
119     you find conflicts among these requirements, you may choose to selectively ignore them
119     .
119     {{difficulty attribute define}}
119
119     ### Difficulty Description:
119 The following are some samples (0 or several).
120 Please ensure that the difficulty level of the samples you generate is harder than these
120     examples.
120 The samples you generate should aim to assess different knowledge and skills compared to
120     the given samples.
120 The format of given samples are not what you should follow.
120 **Please ensure that the sample you create differ substantially from the following samples
120     , so as to maintain diversity in the resulting benchmark.**
120     {{demonstrations}}
120
120     ### Generation Guidelines:
120

```

```

121 **Stage 1: Analyze**
121 In this stage, you should analyze following the steps below and begin with "##Analyses
121 :##". **You need to clearly articulate the analysis content for each step**, which
121 means after completing Stage 1, you should have already produced a question that meets
121 the requirements along with a correct and unique answer.
121 1-1. Analyze the general attributes, difficulty attributes and difficulty description, and
121 think step-by-step about the content needed to construct the question. **Please use
121 your imagination and avoid any obvious overlap with the given samples, either in the
121 specific knowledge points being tested or in the format.**
121 1-2. Start by drafting your question. If you discover any issues with the question or any
121 overlapping parts between the generated question and the given samples during this
121 process, feel free to revise it.
122 1-3. Think through what the correct answer should be. If you discover any issues during
122 this process, repeat the entire Stage 1 process from the beginning.
122 1-4. Identify the plausible and potentially misleading incorrect options that could serve
122 as distractors (at least nine). If you discover any issues during this process, repeat
122 the entire Stage 1 process from the beginning.
122 1-5. Reevaluate your proposed question, answer and options to ensure that: the question
122 meet the given attributes and Difficulty Description (you should compare the generated
122 samples and given samples to verify this); the answer is both correct and unique. If
122 it does not meet these criteria or you are not sure about this, repeat the entire
122 Stage 1 process from the beginning.
122

122 **Stage 2: Generate Sample**
123 In this stage, you should give your generated sample in the right template based on the
123 analyses above.
123 2-1. Generate the question content, begin with "##Question:#".
123 2-2. Generate a step-by-step reasoning process and the corresponding correct answer. Begin
123 with "##Reasoning Path:#". If you find an issue with the question, return to Step
123 2-1 to regenerate the question.
123 2-3. Generate {{OptionNum}} candidates, with only one as the right option. Begin with "##
123 Candidates:#".
123 2-4. Generate the index of the right option, begin with "##Right Option:#".
123

123 ### Output Description:
123 Strictly follow the template below to generate your sample.
124 **Template**
124 ##Analyses:# {{You analyze the provided attributes and outline the process for
124 constructing the question to be generated.}}
124 ##Question:# {{Your generated question content}}
124 ##Reasoning Path:# {{Your step-by-step reasoning process}}
124 ##Candidates:# {{CandidatesDemo}}
124 ##Right Option:# {{Index of the right option, e.g., B}}
124 **Template End**

124

124 Attention: You need to **strictly follow the template** and don't generate any other
125 contents. Begin your response with "##Analyses:#\n1-1. "
125
125

```

1253 **K. Assessment Demands List**

1254 **MATH:**

```

1255 Subset Name: Prealgebra
1256 Assessment Demands:Prealgebra

1257 Subset Name: Algebra
1258 Assessment Demands:Algebra

1259 Subset Name: Number Theory
1260 Assessment Demands:Number Theory
1261
1262
1263
1264

```

```
126 Subset Name: Counting & Probability  
126 Assessment Demands:Counting & Probability  
126  
126 Subset Name: Geometry  
126 Assessment Demands:Geometry  
126  
127 Subset Name: Intermediate Algebra  
127 Assessment Demands:Intermediate Algebra  
127  
127 Subset Name: Precalculus  
127 Assessment Demands:Precalculus  
127
```

MMLU-Pro:

```
127 Subset Name: psychology  
127 Assessment Demands:This benchmark is designed to assess **psychology** abilities while  
127     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
127     **ten multiple-choice questions** as the evaluation format  
128  
128 Subset Name: philosophy  
128 Assessment Demands:This benchmark is designed to assess **philosophy** abilities while  
128     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
128     **ten multiple-choice questions** as the evaluation format  
128  
128 Subset Name: health  
128 Assessment Demands:This benchmark is designed to assess **health** abilities while  
128     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
128     **ten multiple-choice questions** as the evaluation format  
128  
128 Subset Name: history  
128 Assessment Demands:This benchmark is designed to assess **history** abilities while  
128     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
128     **ten multiple-choice questions** as the evaluation format  
129  
129 Subset Name: business  
129 Assessment Demands:This benchmark is designed to assess **business** abilities while  
129     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
129     **ten multiple-choice questions** as the evaluation format  
129  
129 Subset Name: physics  
129 Assessment Demands:This benchmark is designed to assess **physics** abilities while  
129     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
130     **ten multiple-choice questions** as the evaluation format  
130  
130 Subset Name: engineering  
130 Assessment Demands:This benchmark is designed to assess **engineering** abilities while  
130     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
130     **ten multiple-choice questions** as the evaluation format  
130  
130 Subset Name: chemistry  
130 Assessment Demands:This benchmark is designed to assess **chemistry** abilities while  
130     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
130     **ten multiple-choice questions** as the evaluation format  
130  
131 Subset Name: math  
131 Assessment Demands:This benchmark is designed to assess **math** abilities while  
131     simultaneously evaluating knowledge understanding and complex reasoning skills, using  
131     **ten multiple-choice questions** as the evaluation format  
131  
131 Subset Name: computer science  
131 Assessment Demands:This benchmark is designed to assess **computer science** abilities  
131     while simultaneously evaluating knowledge understanding and complex reasoning skills,  
131     using **ten multiple-choice questions** as the evaluation format  
131  
131 Subset Name: biology
```

1320 Assessment Demands: This benchmark is designed to assess **biology** abilities while
1321 simultaneously evaluating knowledge understanding and complex reasoning skills, using
1322 **ten multiple-choice questions** as the evaluation format

1323 Subset Name: economics
1324 Assessment Demands: This benchmark is designed to assess **economics** abilities while
1325 simultaneously evaluating knowledge understanding and complex reasoning skills, using
1326 **ten multiple-choice questions** as the evaluation format

1327 Subset Name: law
1328 Assessment Demands: This benchmark is designed to assess **law** abilities while
1329 simultaneously evaluating knowledge understanding and complex reasoning skills, using
1330 **ten multiple-choice questions** as the evaluation format

1331

1332 **HellaSwag:**

1333

1334 Subset Name: NLI
1335 Assessment Demands: The task is to evaluate the model's commonsense natural language
1336 inference ability. Specifically, each question should present a concrete scenario, and
1337 the model should select the most likely event from the options based on a series of
1338 inferences.

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374