

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

Anonymous ACL submission

## 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, generic and efficient benchmark generators that can construct high-quality benchmarks are widely needed. However, human annotators are constrained by inefficiency, and current LLM-based benchmark generators lack not only generalizability but also a comprehensive evaluation framework for validation and optimization. To fill this gap, we first establish an automated 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. On this basis, we introduce a series of methods to address the identified weaknesses and integrate them as BENCHMARKER. Experiments across multiple LLMs and tasks confirm that BENCHMARKER achieves comparable performance to human-annotated benchmarks on most metrics, highlighting its generalizability and validity. More importantly, it delivers highly consistent evaluation results across 21 LLMs (e.g., 0.969 Pearson correlation against MMLU-Pro on language understanding task), while incurring minimal overhead (e.g., \$0.005 and 0.38 minutes per sample if using GPT-4o mini as generator).

## 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 pace of releasing available LLMs 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 between massive LLM supply and various application needs, the demand for customized benchmarks is

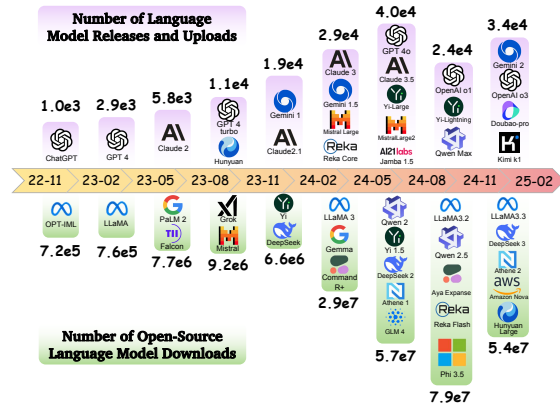


Figure 1: The trends of LLMs released and uploaded, and open-source LLMs downloads per season since the debut of ChatGPT. We obtain the data via the Huggingface API (Appendix D).

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 generalizability across tasks. Meanwhile, the absence of a comprehensive evaluation framework hinders the assessment and optimization of benchmark generators, weakening confidence in LLM-generated benchmarks for real applications. Hence, a comprehensive evaluation framework and a generic LLM-based benchmark generator that can handle any assessment demands and efficiently generate high-quality samples are urgently needed.

To this end, we first establish an automatic evaluation framework with ten criteria for LLM benchmark generators. Notably, we identify the

sources of bias in evaluating label correctness and task relevance using LLM-as-a-judge (Liu et al., 2023) through correlation analysis, and remove them using Multiple Regression. On this basis, we examine the strengths and weaknesses of naively prompting LLM as generic benchmark generator under this evaluation framework. The results reveal that the generated benchmarks exhibit limited lexical and semantic diversity, poor controllability over difficulty, and low label correctness, while showing advantages in high task relevance and behavior diversity. Bearing this in mind, we develop an LLM-based generic benchmark generation method **BENCHMAKER** by integrating existing techniques with newly designed approaches to address the identified issues. Specifically, **BENCHMAKER**: strengthens label correctness 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 from **BENCHMAKER** exhibit better task relevance and difficulty controllability, more challenging question difficulty, and comparable diversity, albeit with slightly lower label correctness. More importantly, they yield highly consistent evaluation results across 21 LLMs (e.g., 0.969 Pearson correlation with MMLU-Pro), while take minimal overhead (e.g., \$0.005 and 0.38 minutes per sample when using GPT-4o mini as generator). We further perform detailed experiments to validate the generalizability 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 label correctness may be imperfect, further enhancing the practicality of **BENCHMAKER** in Appendix A.

## 2 Backgrounds

### 2.1 Data Synthesis

The growth of LLMs’ abilities has led to widespread research on LLM-driven data synthesis,

which demonstrates much better quality and controllability over traditional methods (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 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 framework for benchmark generators. As for implementation, current LLM-based benchmark generators (Perez et al., 2023; Wu et al., 2024; Zhu et al., 2024b; Lei et al., 2023; Li et al., 2025) are constrained by their dependence on seed benchmarks and task specific designs, preventing them from being generic. We construct a comprehensive evaluation framework and develop generic benchmark generation method **BENCHMAKER** to fill these gaps.

### 2.2 Potential Applicable Scenarios of **BENCHMAKER**

Given arbitrary assessment demands  $X$  (e.g., task description, sample type) as input, an automatic generic benchmark generator  $\mathcal{G}$  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, developing **BENCHMAKER** holds significant importance for both scientific research and practical applications within the NLP community.

## 3 Benchmarking Benchmark Generator

While the effectiveness of synthetic training data can be directly evaluated through the trained models, synthetic benchmarks demand a more thorough

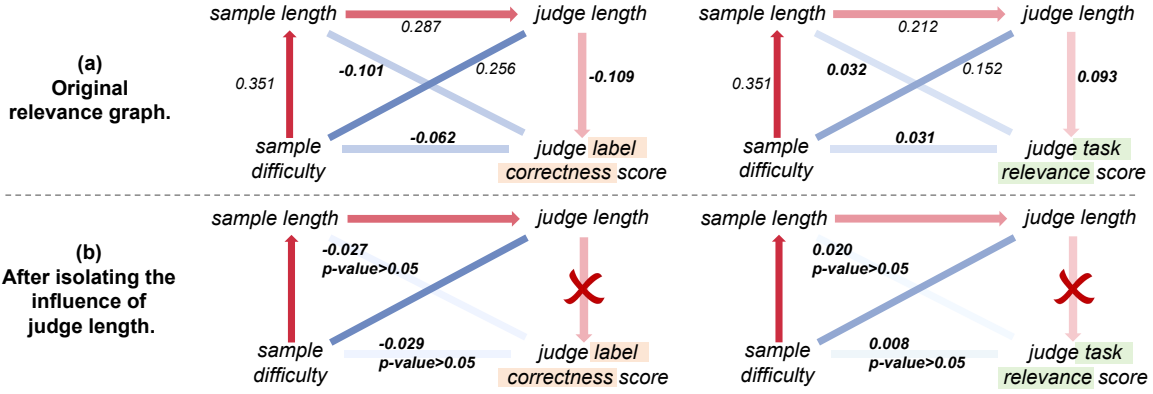


Figure 2: Pearson correlations among key factors of benchmark evaluation and LLM judge scores (label correctness and task relevance, OpenAI o3-mini as the judge). The most relevant path of each subject is highlighted in red.

and multifaceted assessment to ensure their reliability. To this end, we consolidate insights from existing studies to establish a comprehensive ten-criteria evaluation framework for benchmark generators.

### 3.1 Validity

Two key criteria for ensuring the validity of a benchmark are **Label Correctness** and **Task Relevance**. Label correctness measures the accuracy of the generated samples’ label (Wu et al., 2024), while task relevance assesses whether the samples effectively target the capability specified by the assessment demands  $X$  (Perez et al., 2023). For these criteria, previous studies rely on human evaluation (Zhu et al., 2024b) or LLM-as-a-judge (Zheng et al., 2023). However, the former lacks automation, and the latter may introduce bias (Thakur et al., 2024).

To this end, we seek to detect and mitigate possible biases of LLM-as-a-judge that may exist within the framework. We choose OpenAI o3-mini (OpenAI, 2025) as the judge for preliminary experiments with scoring range as  $[0, 1]$  (See prompts for LLM-as-a-judge in Appendix L). Experiments are conducted on the high-quality MATH benchmark (Hendrycks et al., 2021b), where each sample is expected to achieve perfect label correctness and task relevance. Ideally, the scores assigned by the judge should not exhibit any consistency with specific factors. However, as shown in Figure 2-(a), both label correctness and task relevance are significantly correlated ( $p$ -value  $< 0.05$ ) with sample difficulty, sample length, and the length of the judge’s rationale. For each factor, we highlight its weightiest path in red, revealing a possible causal chain: *harder questions lead to longer rationales, requiring judges to conduct lengthier analyses. For label correctness, longer analyses increases the likelihood of judge errors, resulting in lower la-*

*bel correctness ratings. While for task relevance, longer analyses increases the probability of task-relevant words appearing and results in higher task relevance ratings.* To validate the above hypothesis, we control the judge length and respectively calculate the partial correlations (Vallat, 2018) of sample difficulty and sample length with the two criteria. As shown in Figure 2-(b), after isolating the influence of judge length, the effects of other factors are no longer significant ( $p$ -value  $> 0.05$ ). Similar conclusions also hold true when Qwen Plus (Yang et al., 2024) serve as the judge (Figure 7).

Based on the analysis above, the observed biases of the LLM judge are mediated by judgment length. Therefore, for benchmark generators  $\mathcal{G}_{1:|\mathcal{G}|}$  under evaluation, we derive their unbiased judge results with a Multiple Regression model. When calculating label correctness (and similarly task relevance), we first obtain the mean score of  $\mathcal{G}_i$  judged by the LLM and treat it as the dependent variable  $f(i)$ . We then set the generator indices as dummy variables  $\beta_i$  and include judge length as a covariate:

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

Here,  $\beta_{len}$  quantifies the effect of judge length bias. After fitting,  $\beta_i$  represents the average debiased score of  $\mathcal{G}_i$ , which serves as the metric for label correctness.  $\epsilon$  is an offset that adjusts the human-annotated benchmark’s  $\beta_i$  to 1. Thus, the score of  $\mathcal{G}_i$  may exceed 1 if it outperforms the human-annotated counterpart.

### 3.2 Diversity

The diversity of the benchmark determines the extent to which evaluation results can reflect the true model capability across the assessed domain.

**Lexical Diversity** reflects vocabulary richness in benchmarks (Yu et al., 2023). Traditional metrics

like vocabulary size and self-BLEU (Zhu et al., 2018) used in Wu et al. (2024); Yu et al. (2023) are biased by sample length (Guo and Vosoughi, 2023). We use unbiased word frequency entropy (Montahaei et al., 2019) as the metric to evaluate lexical diversity.

**Semantic Diversity** quantifies a benchmark’s semantic comprehensiveness (Chan et al., 2024). We calculate the average Euclidean distance between semantic embeddings of samples as the metric. Specifically, we use text-embedding-ada-002 (OpenAI, 2022) as embedding model for experiments.

**Behavior Diversity** measures whether different samples lead to different model behaviors (Vivek et al., 2024). If a list of evaluated models (denoted as  $\mathcal{M}_{1:|\mathcal{M}|}$ , see Appendix H for detailed list) present the same correctness pattern across two samples (e.g., [ $\checkmark$ ,  $\times$ ,  $\checkmark$ ,  $\checkmark$ ,  $\times$ ] for both), the two samples are essentially interchangeable. We represent the behavior embedding of each sample by the correctness pattern vector across the evaluated models on it. Given the binary nature of this embedding ( $\{0,1\}$ ), we quantify behavior diversity using the average pairwise Hamming distance (Hamming, 1950) among samples. A higher value indicates greater sample irreplaceability.

### 3.3 Difficulty

The difficulty attribute of benchmarks is particularly significant in an era of rapid model evolution.

**Difficulty Controllability** refers to assigning accurate difficulty labels to the samples (e.g., MATH (Hendrycks et al., 2021b)). These labels enable the benchmark to be divided into subsets for more targeted evaluation of models with varying capabilities. For each sample, we use the average error rate of  $\mathcal{M}_{1:|\mathcal{M}|}$  as the ground truth for difficulty label. Based on this, we compute the Spearman correlation between the difficulty labels from the benchmark and the ground truth as the metric for difficulty controllability.

**Difficulty Boundary** denotes the difficulty of the hardest subset of a benchmark (Patel et al., 2025). With the growing strength of LLMs, the performance of the most advanced LLMs on simpler benchmarks has reached saturation (Hendrycks et al., 2021a), making it difficult to differentiate their capabilities. Consequently, more challenging benchmarks (Wang et al., 2024b) are continuously introduced to evaluate the latest LLMs. Thus, we propose assessing the average error rate of  $\mathcal{M}_{1:|\mathcal{M}|}$  on the hardest subset of certain benchmark to mea-

sure its difficulty boundary.

### 3.4 Benchmark-Level Metrics

Lastly, benchmark-level metrics are used to holistically assess the benchmark generation methods.

**Effectiveness.** While the earlier criteria assess benchmark from various aspects, a unified metric is required to measure benchmark effectiveness. Taking high-quality human-annotated benchmark as a reference, we examine if generated benchmark under identical assessment demands can deliver equivalent evaluation results. We calculate the accuracies of  $\mathcal{M}_{1:|\mathcal{M}|}$  on both generated and human benchmarks and calculate their Pearson correlation as the effectiveness metric (Perlitz et al., 2024).

**Robustness.** Under similar inputs, a robust system should produce comparable outputs. Similarly, we expect a robust benchmark generator to produce benchmarks with equivalent evaluation efficacy for similar assessment demands. Thus, we calculate the accuracy of  $\mathcal{M}_{1:|\mathcal{M}|}$  on benchmarks generated under the original assessment demands and that rewritten by an LLM (we use GPT-4o (Hurst et al., 2024) for experiments), 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 generating benchmarks of a certain size.

By establishing this comprehensive evaluation framework, the strengths and weaknesses of benchmark generators can be thoroughly assessed and compared. We further discuss the significance of evaluating benchmark quality from multiple dimensions in Appendix C.

## 4 Development of BenchMaker

In this section, we first analyze the pros and cons of directly prompting the LLM  $\bar{\mathcal{M}}$  as generic benchmark generator (with assessment demands  $X$  as the sole input) in §4.1. On this basis, we refine its weaknesses in the following sections, leading to the development of BENCHMAKER.

### 4.1 Pros and Cons of Directly Prompting

We select three tasks for our experiments, each associated with a high-quality human-annotated benchmark for comparison: Math Reasoning (MATH (Hendrycks et al., 2021b)), Language Understanding (MMLU-Pro (Wang et al., 2024b)),

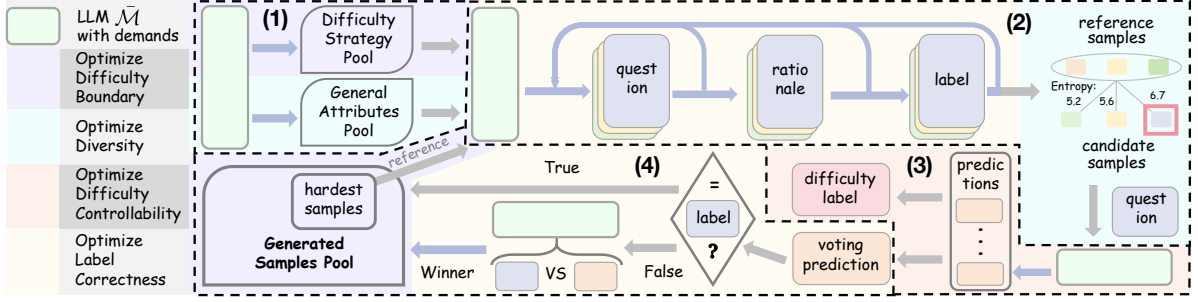


Figure 3: Workflow of BENCHMARKER: (1) Generator model  $\bar{M}$  takes the assessment demands  $X$  to generate an attribute pool and a difficulty strategy pool; (2) Given  $X$ , sampled attributes, difficulty strategy, and hardest samples,  $\bar{M}$  applies *Stepwise Self-correction*, *Difficulty Diffusion Mechanism*, *Difficulty Strategy Guidance*, and *Attribute-based Generation* to generate  $L$  candidate samples, from which the most diverse one is selected via *In-batch Diversity Boosting*; (3) The generated question is input into  $\bar{M}$  to produce  $T$  predictions, based on which  $\gamma$  is computed as the difficulty label; (4) The label is further refined using *Conflict Guided Contrastive Discrimination*.

and Commonsense Reasoning (HellaSwag (Zellers et al., 2019)). Following the sample format of human-annotated benchmarks, the first task adopts a free-form question-answering format, while the latter two take the form of multiple-choice questions. The evaluation capabilities and sample format are conveyed to generator LLM  $\bar{M}$  through the assessment demands  $X$ , as shown in Appendix N. We adopt the prompt<sub>base</sub> in Appendix L to guide  $\bar{M}$  (GPT-4o mini by default) in generating credible and diverse samples  $s_{1:|\mathcal{D}_{human}|}$ :

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

where  $q_i, r_i, a_i$  denote question, rationale, and label, respectively. We proportionally adjust the difficulty level  $l$  in the prompt following the difficulty definition in Hendrycks et al. (2021a) (see descriptions in Appendix G), and select samples with the highest difficulty level to form the hardest subset. As shown in Table 1, compared to  $\mathcal{D}_{human}$ , directly prompting LLM as generic benchmark generator demonstrates poorer label correctness, lower lexical and semantic diversity, weaker difficulty controllability, and less challenging subset. Meanwhile, we also observe its advantages in better task relevance<sup>1</sup>, greater behavior diversity, and improved efficiency.

## 4.2 Label Correctness Optimization

To enhance label correctness, 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

<sup>1</sup>We 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.

following two BenchMaker-compatible techniques to optimize label correctness.

**Stepwise Self-correction.** Since errors might occur at any step during the generation of  $\{q_i, r_i, a_i\}$ , we instruct  $\bar{M}$  to validate the content at each step. If an error is detected,  $\bar{M}$  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  $\bar{M}$  not only act as a judge but also as a test-taker to identify potential errors. Let  $\bar{M}$  predicts  $q_i$   $T$  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  $\bar{M}$  conduct a contrastive discrimination between  $r_i$  and  $\hat{r}_i$  to determine the final rationale and label for  $s_i$ .

## 4.3 Difficulty Optimization

**Difficulty Controllability.** From §4.1, we know that the LLM’s ability to control the difficulty of generated samples is limited. In particular, for the language understanding task, the Spearman correlation between the actual and estimated difficulty of the samples is near-zero. To further explore this, we examine LLM’s difficulty perception by asking  $\bar{M}$  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

407 from the perspective of test-taker:

$$408 \quad \gamma_i = \frac{1}{T} \sum_{j=1}^T \mathbf{1}_{\bar{a}_i^j \neq a_i} \quad (3)$$

409 By taking the inconsistency between  $\bar{\mathcal{M}}$ 's predic-  
410 tions  $\bar{a}_i^{1:T}$  and label  $a_i$  as difficulty label, the corre-  
411 lation increases to 0.415, suggesting that  $\gamma$  is more  
412 reliable for estimating difficulty label.

413 **Difficulty Diffusion Mechanism.** Given that the  
414 LLM has a certain level of difficulty perception,  
415 we iteratively select the more challenging samples  
416 according to  $\gamma$  from the generated ones as difficulty  
417 references, and instruct  $\bar{\mathcal{M}}$  to generate a more dif-  
418 ficult sample. This allows the sample difficulty to  
419 rise continuously through diffusion. The detailed  
420 algorithm is described in Appendix I.

421 **Difficulty Strategy Guidance.** We further con-  
422 sider providing  $\bar{\mathcal{M}}$  with task-specific difficulty-  
423 control strategies. Specifically, we first instruct  $\bar{\mathcal{M}}$   
424 to give varying strategies for generating samples  
425 of specific difficulty levels based on assessment  
426 demands  $X$  (see examples in Appendix M). For  
427 example, difficult samples generally require more  
428 reasoning steps. With Difficulty Diffusion Mech-  
429 anism, we progressively introduce more difficult  
430 difficulty-control strategies to  $\bar{\mathcal{M}}$  to further extend  
431 the difficulty boundary.

#### 432 4.4 Diversity Optimization

433 The optimization of synthetic data diversity has  
434 been widely studied (Wang et al., 2024a). We con-  
435 duct extensive tests and select the most generic  
436 and effective **Attribute-based Generation Tech-**  
437 **nique** introduced in AttrPrompt (Yu et al., 2023)  
438 for BENCHMARKER. This method explicitly en-  
439 hances the lexical and semantic diversity of bench-  
440 marks by randomly assigning pre-generated (at-  
441 tribute, value) pairs as part of the input for each  
442 sample. Furthermore, we notice that the introduc-  
443 tion of treating the generated samples as difficulty  
444 references might cause sample homogeneity. To  
445 mitigate this, we propose an **In-batch Diversity**  
446 **Boosting** method, where  $\bar{\mathcal{M}}$  generates  $L$  (We set  $L$   
447 as 5 for our default setting) candidate samples and  
448 selects the one with the greatest word frequency en-  
449 tropy difference from the input reference samples.

## 450 5 Experiments and Analyses

451 We run extensive experiments to evaluate BENCH-  
452 MAKER under the proposed evaluation framework.

453 **Settings.** We select the widely used human-  
454 annotated MATH (Hendrycks et al., 2021b) (math-  
455 ematical reasoning), MMLU-Pro (multi-task lan-  
456 guage understanding) (Wang et al., 2024b) and Hel-  
457 laSwag (commonsense reasoning) (Zellers et al.,  
458 2019) as high-quality baseline benchmarks. For the  
459 7 subsets of MATH, the 13 subsets of MMLU-Pro<sup>2</sup>  
460 and HellaSwag, we write simple assessment de-  
461 mands  $X$  respectively (see details in Appendix N)  
462 as inputs for the generator model  $\bar{\mathcal{M}}$ . For each de-  
463 mand, we generate  $N$  (default as 500) samples and  
464 randomly downsample the human-annotated bench-  
465 mark to match the number of generated samples  
466 for fair comparison (when calculating effectiveness,  
467 we avoid downsampling to ensure more accurate  
468 results). Each experiment is repeated three times,  
469 and the average results are reported. We exper-  
470 iment with GPT-4o mini (Hurst et al., 2024) as  
471 the default  $\bar{\mathcal{M}}$  and also explore the performance of  
472 GPT-4o and Claude 3.5 Haiku (Anthropic, 2024).  
473 The decoding temperature is set to 1. To mitigate  
474 the self-enhancement bias (Zheng et al., 2023) of  
475 LLM-as-a-judge, we substitute the generators with  
476 OpenAI o3-mini (OpenAI, 2025) as the judge. Fol-  
477 lowing the guidelines of Perlitiz et al. (2024), we  
478 use 21 LLMs listed in Appendix H as models under  
479 evaluation. We choose the direct prompt method  
480 introduced in §4.1 and the widely used AttrPrompt  
481 (Yu et al., 2023) method as baselines.

### 482 5.1 Comparison with Human-annotated 483 Benchmark

484 Overall (Table 1), while benchmarks from BENCH-  
485 MAKER exhibit slightly reduced label correctness,  
486 they perform on par with human-annotated bench-  
487 marks in lexical and semantic diversity, and surpass  
488 them in task relevance, behavior diversity, difficulty  
489 controllability, and efficiency.

490 **Effectiveness.** The primary goal of benchmark-  
491 ing is to assign accurate scores to models under  
492 evaluation, facilitating capability differentiation.  
493 The benchmarking results of BENCHMARKER align  
494 closely with human-annotated benchmarks, with an  
495 average of 0.955 linear correlation (Pearson) and a  
496 remarkable 0.967 for rank-order correlation (Spear-  
497 man), highlighting its outstanding effectiveness.

498 **Robustness.** Under evaluation demands where  
499 semantic equivalence is maintained but linguistic  
500 styles vary, the benchmarks exhibit nearly identical  
501 assessment efficacy, with an average Pearson corre-

<sup>2</sup>Excluding the type 'other'.

Methods	Label Correct	TaskRelevance	LexicalDiv	SemanticDiv	BehaviorDiv	DifControl	DifBoundary	Effectiveness	Robustness	Efficiency
	Unbias Score↑	Unbias Score↑	Entropy↑	EuclideanDis↑	Hamm-ingDis↑	Spea-rman↑	Error Rate↑	Pear-son↑	Pear-son↑	\$/item, min/item↓
Math Reasoning										
Human Benchmark	<b>1.000</b>	1.000	8.032	0.668	0.363	0.178	0.652	-	-	high
AttrPrompt	0.638	1.145	8.292	0.672	0.364	0.141	0.599	0.747	0.981	0.002, 0.19
Direct Prompt	0.692	1.120	7.105	0.621	0.366	0.115	0.570	0.654	0.989	0.002, 0.17
+InBatchDivBoost	0.660	1.097	8.627	0.676	0.365	0.172	0.548	0.780	0.982	0.003, 0.20
+StepSelfCorrect	0.928	1.098	8.640	0.675	0.377	0.170	0.481	0.792	0.984	0.003, 0.23
+ConflictConDisc	0.987	1.134	8.671	0.678	0.372	0.202	0.434	0.851	0.987	0.004, 0.36
+DiffControl	0.987	1.134	8.671	0.678	0.372	0.424	0.434	0.851	0.987	0.004, 0.36
+DiffDiffusion	0.971	1.159	8.694	0.678	0.390	0.468	0.601	0.895	0.990	0.005, 0.39
BenchMaker	0.922	1.151	<b>8.930</b>	<b>0.680</b>	0.407	<b>0.450</b>	0.650	0.931	<b>0.988</b>	<b>0.005, 0.42</b>
BenchMaker <sub>4o</sub>	0.931	<b>1.174</b>	8.866	0.675	0.389	0.444	0.661	<b>0.933</b>	0.983	0.082, 1.10
BenchMaker <sub>haiku</sub>	0.875	1.052	8.895	0.677	<b>0.412</b>	0.386	<b>0.684</b>	0.884	0.971	0.026, 0.56
Language Understanding										
Human Benchmark	1.000	1.000	<b>10.404</b>	<b>0.731</b>	0.311	0.000	<b>0.654</b>	-	-	high
AttrPrompt	0.867	1.185	9.916	0.735	0.390	0.079	0.561	0.890	0.988	0.002, 0.17
Direct Prompt	0.882	1.174	9.608	0.726	0.394	0.037	0.532	0.867	0.989	0.002, 0.16
BenchMaker	<b>1.032</b>	<b>1.211</b>	10.166	0.728	<b>0.397</b>	<b>0.461</b>	0.642	<b>0.969</b>	0.982	<b>0.005, 0.38</b>
Commonsense Reasoning										
Human Benchmark	<b>1.000</b>	1.000	<b>9.167</b>	0.655	0.371	0.000	0.481	-	-	high
AttrPrompt	0.884	1.064	8.763	0.658	0.392	0.121	0.580	0.858	0.974	0.002, 0.19
Direct Prompt	0.878	1.057	8.165	0.660	0.384	0.062	0.536	0.844	0.979	0.002, 0.17
BenchMaker	0.987	<b>1.084</b>	9.052	<b>0.663</b>	<b>0.404</b>	<b>0.445</b>	<b>0.634</b>	<b>0.946</b>	<b>0.984</b>	<b>0.005, 0.40</b>

Table 1: Overall results under the proposed evaluation framework. For each setting, We run 3 times and report the average results. GPT-4o mini serves as the default generator. Values in bold denote the best results between Human Benchmark and BENCHMARKER. The blue lines are results of baselines and the yellow lines are ablation studies.

502 lation of 0.987. This demonstrates the robustness  
503 of BENCHMARKER to diverse inputs and ensures  
504 that users with different linguistic preferences can  
505 obtain consistent evaluation results.

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

513 **Generalizability.** Experimental results demon-  
514 strate that BENCHMARKER exhibits strong general-  
515 ization across task types and generators, achieving  
516 an effectiveness of no less than 0.884 (Claude 3.5  
517 Haiku). Compared to GPT-4o, GPT-4o mini proves  
518 to be a more cost-effective benchmark generator.

## 5.2 Ablation Studies

519 We validate the effectiveness of different tech-  
520 niques by sequentially integrating them to the Di-  
521 rect Prompt baseline. Since the number of ablation  
522 settings is large and the computational cost per set-  
523 ting is substantial, we choose the Math Reasoning  
524

task for ablation studies, as shown in Table 1.

525 **Diversity.** Compared to Direct Prompt, both  
526 Attribute-based Generation Technique and In-batch  
527 Diversity Boosting enhance lexical and semantic  
528 diversity. Noticeably, behavior diversity stays con-  
529 stant, suggesting that surface-level variation does  
530 not inherently result in a sample distribution more  
531 capable of distinguishing between model behaviors.  
532 Meanwhile, the diversity improvement leads to a  
533 slight drop in label correctness, possibly because  
534 of the attributes constraints.

535 **Label Correctness.** After applying Stepwise Self-  
536 correction and Conflict Guided Contrastive Dis-  
537 crimination, we observe a sustained improvement  
538 in label correctness. Meanwhile, we notice that  
539 the difficulty of the hardest subset decreases from  
540 0.548 to 0.434. This may be due to the fact that a  
541 high label error rate often results in performance  
542 being underestimated. Consequently, once the la-  
543 bels are corrected, the accuracy can better reflect  
544 the actual difficulty of the benchmark.

545 **Difficulty Controllability.** By treating the gen-  
546 erator as the test-taker and using its error rate as  
547 the difficulty label, we achieve more precise con-  
548

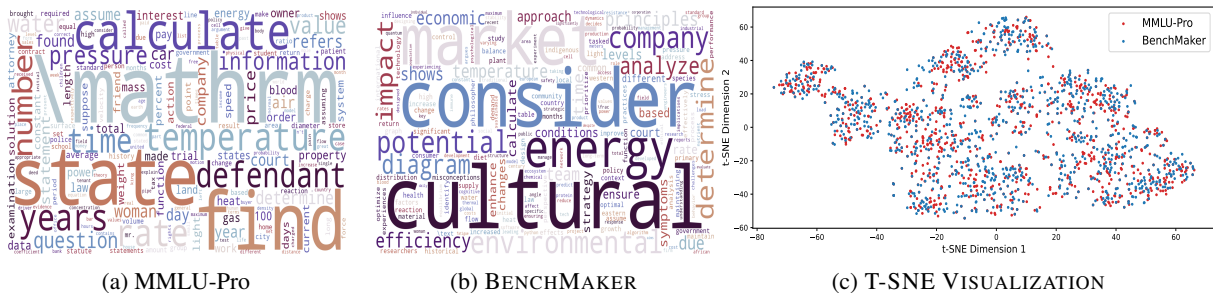


Figure 4: (a)-(b): Word cloud of MMLU-Pro and the benchmark generated by BENCHMARKER under similar assessment demands. (c): T-SNE results on the text embeddings of both benchmarks.

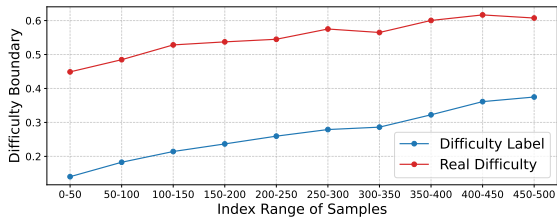


Figure 5: Trends of real and annotated difficulty over the index.

control over sample difficulty (Spearman correlation of 0.424). Considering the observed weak difficulty perception of LLMs, we hypothesize that this improvement stems from the role shift, which requires the model to engage in explicit reasoning, along with the adoption of prediction-label inconsistency as an objective metric.

**Difficulty Boundary.** With the proposed Difficulty Diffusion Mechanism and Difficulty Strategy Guidance, the difficulty boundary is significantly extended, as evidenced by an increase in error rate from 0.434 to 0.650. Additionally, we examine how actual difficulty and difficulty label change with generation order. As shown in Figure 5, both increase steadily, confirming the effectiveness of the Difficulty Diffusion Mechanism.

### 5.3 A Closer Look at Generated Benchmark

After metric analysis, we perform a more thorough examination of BENCHMARKER. Some of the generated samples are shown in Appendix J and the model performance on generated benchmarks are shown in Appendix E.

**Lexical and Semantic.** First, despite the obvious differences in word distribution between the generated benchmark and MMLU-Pro (Figure 4), it remains closely aligned with the domains covered by MMLU-Pro, demonstrating strong task relevance. Meanwhile, the semantic alignment between the two is more pronounced (Figure 4c). Notably, the input demands (Appendix N) do not mention any information related to MMLU-Pro, effectively pre-

Sample Number $N$	125	250	375	500
Effectiveness	0.898	0.940	0.961	0.969

Table 2: The impact of dataset size  $N$ .

venting the model from achieving a high degree of alignment by memorizing and replicating samples from MMLU-Pro.

**Actual Label Correctness.** While the detected potential bias of LLMs in judging label correctness has been addressed, we still carry out a manual review on 80 randomly chosen samples. We find that 3 samples have incorrect labels, 3 samples' questions are confusing, 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) BENCHMARKER still has room for improvement in label correctness; (2) LLM-as-a-judge can serve as a partial proxy for human evaluation.

**Benchmark Size Analyses.** We investigate the impact of dataset size  $N$  on Language Understanding task. As shown in Table 2, with the increases in  $N$ , the improvement in effectiveness gradually slows down, and acceptable effectiveness is already achieved at relatively small values of  $N$ . Therefore, in practical applications, we recommend selecting an appropriate  $N$  based on budget constraints.

## Conclusions

The rapid evolution of LLMs has driven an urgent demand for an automated generic benchmark generator. To this end, we propose a comprehensive framework for benchmark generator evaluation, based on which we develop BENCHMARKER method for reliable, generic, and efficient benchmark generation. Experiments across multiple tasks and LLMs demonstrate that BENCHMARKER achieves human-aligned benchmark quality, with superior efficiency and generalizability.

## 614 Limitations

615 Although we have verified the generalizability of  
616 BenchMaker across multiple settings (3 tasks, 3  
617 generator LLMs, 21 tester LLMs), it still does  
618 not cover all scenarios. We anticipate that Bench-  
619 Maker can undergo more extensive evaluation in  
620 real-world settings, which will help us further opti-  
621 mize it in the future.

622 We emphasize that the primary goal of this work  
623 is to systematically examine the current landscape  
624 of benchmark synthesis with LLMs and to pro-  
625 vide initial mitigation strategies for the key issues  
626 observed in practice. We believe that designing  
627 a truly reliable and generic framework for LLM-  
628 based benchmark synthesis will require sustained  
629 efforts from the broader community. This includes,  
630 but is not limited to, (i) achieving better evaluation-  
631 goal alignment through iterative system–human  
632 interaction, (ii) identifying, quantifying, and miti-  
633 gating potential biases introduced during synthe-  
634 sis and evaluation, and (iii) striking a principled  
635 balance among controllability, cost, and scalabil-  
636 ity. These directions are beyond the scope of the  
637 present study. Nevertheless, we view BenchMaker  
638 as a foundational step that offers reusable insights  
639 and practical guidance, and we hope it can help  
640 catalyze further advances in this emerging area.

## 641 Ethics Statement

642 All of the datasets used in this study were publicly  
643 available. We confirm that the datasets we used did  
644 not contain any harmful content and was consistent  
645 with their intended use (research). We have cited  
646 the datasets and relevant works used in this study.  
647 During the research process, AI was used solely  
648 to assist with code development and to help identi-  
649 fy and correct writing errors in the manuscript;  
650 no AI usage prohibited by relevant policies was  
651 employed.

## 652 References

- 653 Anthropic. 2024. Claude 3.5.  
654 [https://www.anthropic.com/news/  
655 3-5-models-and-computer-use](https://www.anthropic.com/news/3-5-models-and-computer-use).
- 656 Simone Balloccu, Patrícia Schmidová, Mateusz Lango,  
657 and Ondrej Dusek. 2024. **Leak, cheat, repeat: Data  
658 contamination and evaluation malpractices in closed-  
659 source llms**. In *Proceedings of the 18th Conference  
660 of the European Chapter of the Association for Com-  
661 putational Linguistics, EACL 2024 - Volume 1: Long*

*Papers, St. Julian's, Malta, March 17-22, 2024*, pages  
662 67–93. Association for Computational Linguistics. 663

Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and  
664 Dong Yu. 2024. **Scaling synthetic data creation with  
665 1,000,000,000 personas**. *CoRR*, abs/2406.20094. 666

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,  
667 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,  
668 Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang,  
669 Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie.  
670 2024. **A survey on evaluation of large language mod-  
671 els**. *ACM Trans. Intell. Syst. Technol.*, 15(3):39:1–  
672 39:45. 673

Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego  
674 Chicharro, Evan Chen, Alex Gunning, Caroline Falk-  
675 man Olsson, Jean-Stanislas Denain, Anson Ho,  
676 Emily de Oliveira Santos, Olli Järviemi, Matthew  
677 Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla,  
678 Qiuyu Ren, Elizabeth Pratt, Lionel Levine, Grant  
679 Barkley, and 5 others. 2024. **Frontiermath: A bench-  
680 mark for evaluating advanced mathematical reason-  
681 ing in AI**. *CoRR*, abs/2411.04872. 682

Xiaobo Guo and Soroush Vosoughi. 2023. **Length does  
683 matter: Summary length can bias summarization met-  
684 rics**. In *Proceedings of the 2023 Conference on  
685 Empirical Methods in Natural Language Process-  
686 ing, EMNLP 2023, Singapore, December 6-10, 2023*,  
687 pages 15869–15879. Association for Computational  
688 Linguistics. 689

Richard W Hamming. 1950. Error detecting and error  
690 correcting codes. *The Bell system technical journal*,  
691 29(2):147–160. 692

Dan Hendrycks, Collin Burns, Steven Basart, Andy  
693 Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-  
694 hardt. 2021a. **Measuring massive multitask language  
695 understanding**. In *9th International Conference on  
696 Learning Representations, ICLR 2021, Virtual Event,  
697 Austria, May 3-7, 2021*. OpenReview.net. 698

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul  
699 Arora, Steven Basart, Eric Tang, Dawn Song, and  
700 Jacob Steinhardt. 2021b. **Measuring mathematical  
701 problem solving with the MATH dataset**. In *Pro-  
702 ceedings of the Neural Information Processing Sys-  
703 tems Track on Datasets and Benchmarks 1, NeurIPS  
704 Datasets and Benchmarks 2021, December 2021, vir-  
705 tual*. 706

Jie Huang, Xinyun Chen, Swaroop Mishra,  
707 Huaixiu Steven Zheng, Adams Wei Yu, Xinyun  
708 Song, and Denny Zhou. 2024. **Large language  
709 models cannot self-correct reasoning yet**. In *The  
710 Twelfth International Conference on Learning  
711 Representations, ICLR 2024, Vienna, Austria, May  
712 7-11, 2024*. OpenReview.net. 713

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam  
714 Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,  
715 Akila Welihinda, Alan Hayes, Alec Radford, and 1  
716 others. 2024. **Gpt-4o system card**. *arXiv preprint  
717 arXiv:2410.21276*. 718

719	Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023. <a href="#">Towards mitigating LLM hallucination via self reflection</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1827–1843. Association for Computational Linguistics.	775
720		776
721		777
722		778
723		779
724		
725		
726	Fangyu Lei, Qian Liu, Yiming Huang, Shizhu He, Jun Zhao, and Kang Liu. 2023. <a href="#">S3eval: A synthetic, scalable, systematic evaluation suite for large language models</a> . <i>CoRR</i> , abs/2310.15147.	
727		
728		
729		
730	Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. <a href="#">Retrieval-augmented generation for knowledge-intensive NLP tasks</a> . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	
731		
732		
733		
734		
735		
736		
737		
738		
739	Jiatong Li, Renjun Hu, Kunzhe Huang, Yan Zhuang, Qi Liu, Mengxiao Zhu, Xing Shi, and Wei Lin. 2024a. <a href="#">Perteval: Unveiling real knowledge capacity of llms with knowledge-invariant perturbations</a> . <i>CoRR</i> , abs/2405.19740.	
740		
741		
742		
743		
744	Sihang Li, Jin Huang, Jiayi Zhuang, Yaorui Shi, Xiaochen Cai, Mingjun Xu, Xiang Wang, Linfeng Zhang, Guolin Ke, and Hengxing Cai. 2024b. <a href="#">Scilitlm: How to adapt llms for scientific literature understanding</a> . <i>CoRR</i> , abs/2408.15545.	
745		
746		
747		
748		
749	Xiang Lisa Li, Farzaan Kaiyom, Evan Zheran Liu, Yifan Mai, Percy Liang, and Tatsunori Hashimoto. 2025. <a href="#">Autobench: Towards declarative benchmark construction</a> . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	
750		
751		
752		
753		
754		
755	Zenan Li, Zhi Zhou, Yuan Yao, Yu-Feng Li, Chun Cao, Fan Yang, Xian Zhang, and Xiaoxing Ma. 2024c. <a href="#">Neuro-symbolic data generation for math reasoning</a> . <i>CoRR</i> , abs/2412.04857.	
756		
757		
758		
759	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen, and Chenguang Zhu. 2023. <a href="#">G-eval: NLG evaluation using gpt-4 with better human alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 2511–2522. Association for Computational Linguistics.	
760		
761		
762		
763		
764		
765		
766		
767	Lin Long, Rui Wang, Ruixuan Xiao, Junbo Zhao, Xiao Ding, Gang Chen, and Haobo Wang. 2024. <a href="#">On llms-driven synthetic data generation, curation, and evaluation: A survey</a> . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 11065–11082. Association for Computational Linguistics.	
768		
769		
770		
771		
772		
773		
774		
	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Qingwei Lin, Jianguang Lou, Shifeng Chen, Yansong Tang, and Weizhu Chen. 2024a. <a href="#">Arena learning: Build data flywheel for llms post-training via simulated chatbot arena</a> . <i>CoRR</i> , abs/2407.10627.	
	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2024b. <a href="#">Wizardcoder: Empowering code large language models with evolve-instruct</a> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	780
		781
		782
		783
		784
		785
		786
	Gaurav Maheshwari, Dmitry Ivanov, and Kevin El Haddad. 2024. <a href="#">Efficacy of synthetic data as a benchmark</a> . <i>CoRR</i> , abs/2409.11968.	787
		788
		789
	Ehsan Montahaei, Danial Alihosseini, and Mahdieh Soleymani Baghshah. 2019. <a href="#">Jointly measuring diversity and quality in text generation models</a> . <i>CoRR</i> , abs/1904.03971.	790
		791
		792
		793
	OpenAI. 2022. <a href="#">text-embedding-ada-002</a> . <a href="https://platform.openai.com/docs/guides/embeddings">https://platform.openai.com/docs/guides/embeddings</a> .	794
		795
		796
	OpenAI. 2025. <a href="#">Openai o3 and o4-mini system card</a> .	797
	Arkil Patel, Siva Reddy, and Dzmitry Bahdanau. 2025. <a href="#">How to get your llm to generate challenging problems for evaluation</a> . <i>arXiv preprint arXiv:2502.14678</i> .	798
		799
		800
	Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, and 44 others. 2023. <a href="#">Discovering language model behaviors with model-written evaluations</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13387–13434. Association for Computational Linguistics.	801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
	Yotam Perlitz, Ariel Gera, Ofir Arviv, Asaf Yehudai, Elron Bandel, Eyal Shnarch, Michal Shmueli-Scheuer, and Leshem Choshen. 2024. <a href="#">Benchmark agreement testing done right: A guide for LLM benchmark evaluation</a> . <i>CoRR</i> , abs/2407.13696.	813
		814
		815
		816
		817
	Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. 2024. <a href="#">A survey on self-evolution of large language models</a> . <i>CoRR</i> , abs/2404.14387.	818
		819
		820
		821
		822
	Aman Singh Thakur, Kartik Choudhary, Venkat Srinik Ramayapally, Sankaran Vaidyanathan, and Dieuwke Hupkes. 2024. <a href="#">Judging the judges: Evaluating alignment and vulnerabilities in llms-as-judges</a> . <i>CoRR</i> , abs/2406.12624.	823
		824
		825
		826
		827
	Raphael Vallat. 2018. <a href="#">Pingouin: statistics in python</a> . <i>J. Open Source Softw.</i> , 3(31):1026.	828
		829

830	Rajan Vivek, Kawin Ethayarajh, Diyi Yang, and Douwe Kiela. 2024. <a href="#">Anchor points: Benchmarking models with much fewer examples</a> . In <i>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</i> , pages 1576–1601. Association for Computational Linguistics.	
838	Ke Wang, Jiahui Zhu, Minjie Ren, Zeming Liu, Shiwei Li, Zongye Zhang, Chenkai Zhang, Xiaoyu Wu, Qiqi Zhan, Qingjie Liu, and Yunhong Wang. 2024a. <a href="#">A survey on data synthesis and augmentation for large language models</a> . <i>CoRR</i> , abs/2410.12896.	
843	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023a. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
850	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. <a href="#">Self-instruct: Aligning language models with self-generated instructions</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13484–13508. Association for Computational Linguistics.	
859	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024b. <a href="#">Mmlu-pro: A more robust and challenging multi-task language understanding benchmark</a> . <i>CoRR</i> , abs/2406.01574.	
866	Siyuan Wu, Yue Huang, Chujie Gao, Dongping Chen, Qihui Zhang, Yao Wan, Tianyi Zhou, Xiangliang Zhang, Jianfeng Gao, Chaowei Xiao, and Lichao Sun. 2024. <a href="#">Unigen: A unified framework for textual dataset generation using large language models</a> . <i>CoRR</i> , abs/2406.18966.	
872	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, and 43 others. 2024. <a href="#">Qwen2 technical report</a> . <i>CoRR</i> , abs/2407.10671.	
879	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024. <a href="#">Meta-math: Bootstrap your own mathematical questions for large language models</a> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	
	Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander J. Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. 2023. <a href="#">Large language model as attributed training data generator: A tale of diversity and bias</a> . In <i>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</i> .	887 888 889 890 891 892 893 894
	Peiwen Yuan, Yueqi Zhang, Shaoxiong Feng, Yiwei Li, Xinglin Wang, Jiayi Shi, Chuyi Tan, Boyuan Pan, Yao Hu, and Kan Li. 2025. <a href="#">Beyond one-size-fits-all: Tailored benchmarks for efficient evaluation</a> . <i>arXiv preprint arXiv:2502.13576</i> .	895 896 897 898 899
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. <a href="#">Hellaswag: Can a machine really finish your sentence?</a> In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 4791–4800. Association for Computational Linguistics.	900 901 902 903 904 905 906 907
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. <a href="#">Judging llm-as-a-judge with mt-bench and chatbot arena</a> . In <i>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</i> .	908 909 910 911 912 913 914 915 916
	Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. 2024a. <a href="#">Dyval: Dynamic evaluation of large language models for reasoning tasks</a> . In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . OpenReview.net.	917 918 919 920 921 922
	Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. 2024b. <a href="#">Dyval 2: Dynamic evaluation of large language models by meta probing agents</a> . <i>CoRR</i> , abs/2402.14865.	923 924 925 926
	Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. <a href="#">Texygen: A benchmarking platform for text generation models</a> . In <i>The 41st International ACM SIGIR Conference on Research &amp; Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018</i> , pages 1097–1100. ACM.	927 928 929 930 931 932 933
	<b>A Results Credibility under Noised Labels.</b>	934 935
	Since label correctness of the generated benchmark cannot be ensured, we are curious about the effects of these noised samples: Assume a generated benchmark of size $N$ contains a fraction $K$ of randomly labeled samples. Given the presence of label noise, we suppose the probability that both models $A$ and $B$ agree with the provided labels on the	936 937 938 939 940 941 942

corrupted samples is  $p$ , and their true accuracies on clean samples are  $\theta_A$  and  $\theta_B$ . The expected observed accuracies satisfy

$$\pi_A = (1 - K)\theta_A + Kp, \quad \pi_B = (1 - K)\theta_B + Kp, \quad (4)$$

For large  $N$  we test  $H_0 : \theta_A \leq \theta_B$  versus  $H_1 : \theta_A > \theta_B$  with the one-sided  $z$ -statistic (see details in Appendix A.1)

$$\begin{aligned} z &= \frac{\hat{\theta}_A - \hat{\theta}_B}{\sqrt{[\hat{\theta}_A(1 - \hat{\theta}_A) + \hat{\theta}_B(1 - \hat{\theta}_B)]/N}} \\ &= \frac{\frac{\hat{\theta} = \frac{\hat{\pi} - Kp}{1 - K}}{(1 - K)(\hat{\pi}_A - \hat{\pi}_B)} \sqrt{(1 - K)^2[\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B)]/N}}{\sqrt{(\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B))/N}} \stackrel{H_0}{\sim} \mathcal{N}(0, 1) \end{aligned} \quad (5)$$

**This means that we can directly compute  $z$  using the observed accuracies  $\hat{\pi}_A$  and  $\hat{\pi}_B$ , based on which the corresponding  $p$ -value can be calculated to determine whether the conclusion that model  $A$  outperforms model  $B$  ( $H_1$ ) is credible.**

We also observe that the factor  $K$  vanishes completely in the computation of  $z$ , implying that a certain proportion of label noise in the benchmark will not affect the statistical significance of model ranking.

### A.1 Proof of the one-sided $z$ -statistic in Eq. (5)

Let  $N$  be the benchmark size,  $K \in [0, 1]$  the fraction of randomly labelled items, suppose that the probability that both models hit the corrupted items is  $p$ ,  $\theta_A, \theta_B$  the true accuracies on the clean subset, and  $\hat{\pi}_A, \hat{\pi}_B$  the empirically observed accuracies on the noisy benchmark. Throughout, hats denote sample estimates and subscripts  $A, B$  identify the two models.

#### step 1. Linking noisy and true accuracies

Since only a proportion  $1 - K$  of the labels are correct, the expectations of the observed accuracies satisfy

$$\pi_A = (1 - K)\theta_A + Kp, \quad \pi_B = (1 - K)\theta_B + Kp. \quad (6)$$

Solving for the latent accuracies yields

$$\theta_A = \frac{\pi_A - Kp}{1 - K}, \quad \theta_B = \frac{\pi_B - Kp}{1 - K}. \quad (7)$$

#### step 2. Difference of true accuracies

Subtracting the two equations in (6) eliminates  $p$  and gives

$$\theta_A - \theta_B = \frac{\pi_A - \pi_B}{1 - K}.$$

Hence  $\theta_A > \theta_B$  is equivalent to  $\pi_A > \pi_B$ ; the noise ratio  $K$  plays no role in the sign of the difference.

#### step 3. Sampling distribution of $\hat{\pi}_A - \hat{\pi}_B$

For large  $N$  each  $\hat{\pi}$  is approximately normal by the Central Limit Theorem:

$$\hat{\pi}_A \sim \mathcal{N}\left(\pi_A, \frac{\pi_A(1 - \pi_A)}{N}\right), \quad \hat{\pi}_B \sim \mathcal{N}\left(\pi_B, \frac{\pi_B(1 - \pi_B)}{N}\right),$$

and they are asymptotically independent. Therefore

$$\hat{\pi}_A - \hat{\pi}_B \sim \mathcal{N}\left(\pi_A - \pi_B, \frac{\pi_A(1 - \pi_A) + \pi_B(1 - \pi_B)}{N}\right).$$

#### step 4. One-sided test $H_0 : \theta_A \leq \theta_B$

vs.  $H_1 : \theta_A > \theta_B$

Using  $\hat{\pi}_A - \hat{\pi}_B$  as the test statistic and plugging the estimated means into the variance gives the empirical  $z$ -score reported in Eq. (5):

$$\begin{aligned} z &= \frac{\hat{\pi}_A - \hat{\pi}_B}{\sqrt{[\hat{\pi}_A(1 - \hat{\pi}_A) + \hat{\pi}_B(1 - \hat{\pi}_B)]/N}} \\ &\stackrel{H_0}{\sim} \mathcal{N}(0, 1). \end{aligned}$$

Crucially, all factors of  $(1 - K)$  cancel; the presence of an arbitrary proportion  $K$  of noisy labels does not alter the null distribution of  $z$ . The  $p$ -value of the one-sided test is therefore

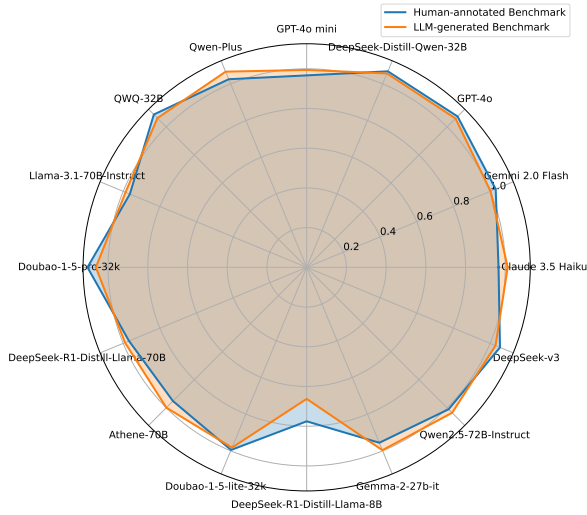
$$p\text{-value} = 1 - \Phi(z),$$

where  $\Phi$  is the standard normal CDF. A small  $p$ -value allows us to reject  $H_0$  and conclude that model  $A$  outperforms model  $B$  on the underlying clean benchmark.

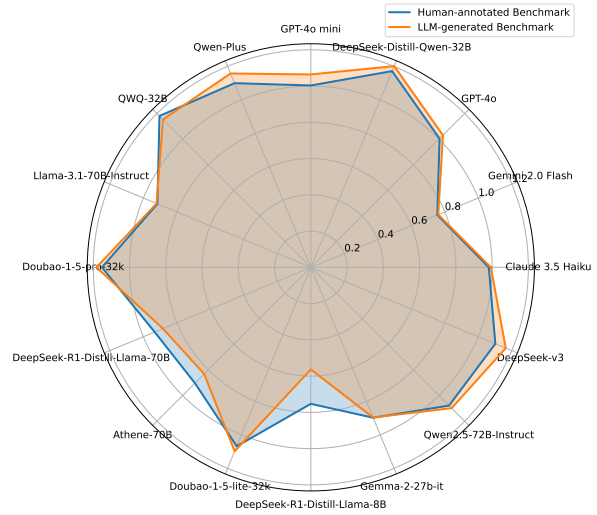
## B Unsuccessful Attempts for Optimizing Benchmark Generator

### B.1 Label Correctness

We explored the widely studied self-correction strategy to improve the label correctness 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 language understanding and instead introduces additional computational overhead.



(a) Language Understanding



(b) Math Reasoning

Figure 6: LLM Performance on human-annotated and LLM-generated benchmarks.

## 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 Attr-Prompt, 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 Significance of a Comprehensive Evaluation Framework for Benchmark Generator

Different scenarios call for different requirements, so it is essential to assess benchmark quality from multiple dimensions.

- When developing benchmark generation methods, tasks with high-quality human-annotated datasets can serve as testbeds, allowing us to directly evaluate the quality of generated benchmarks via effectiveness. Efficiency and robustness are also critical, as they reflect the implementation efficiency and generation stability of different methods.
- When selecting appropriate models for customized scenarios without access to high-quality human annotations, label correctness, task relevance, and diversity become more important. These metrics allow users to quantify whether the generated benchmark accurately evaluates model capabilities under the desired conditions and whether it provides sufficient domain coverage.
- For evaluating state-of-the-art models, the concept of difficulty boundary helps determine whether a benchmark is sufficiently challenging and capable of distinguishing between models, thus indicating whether it has become saturated.

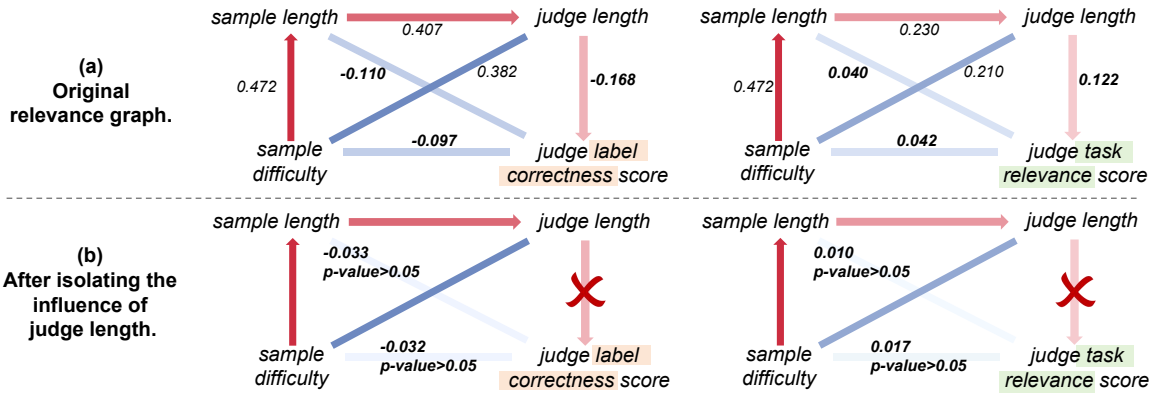


Figure 7: Pearson correlations among key factors of benchmark evaluation and LLM (GPT-4o mini) judge scores (faithfulness and alignment). The most relevant path of each subject is highlighted in red to show the possible causal chain.

- In benchmark compression and efficient evaluation settings, difficulty controllability reflects the accuracy of the provided difficulty labels, while behavior diversity indicates sample irreplaceability. Accurate difficulty labels enable strategies such as uniformly downsampling based on difficulty or difficulty-based sampling tailored to model capabilities (Yuan et al., 2025), leading to more efficient evaluations.

## D 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.

## E Model Performance on Generated Benchmarks

We present the performance of some mainstream LLMs on human-annotated and LLM-generated benchmarks to compare them from the evaluation effectiveness perspective, as shown in Figure 6. We normalize the accuracy of the evaluated models to have a mean of 1, facilitating comparison. As shown, despite some differences, the model-generated benchmark and the human-annotated benchmark yield aligned overall perfor-

mance trends, demonstrating strong evaluation effectiveness.

## F Human Evaluation Settings

Multiple authors jointly conducted the manual check for the Actual Error Rate section, and no extra annotators were employed for our data collection.

## G Difficulty Levels

In Hendrycks et al. (2021a), the questions are categorized into the following four difficulty levels.

- **Elementary Level:** Basic grade-school questions
- **High School Level:** More challenging high-school curriculum questions
- **College Level:** Undergraduate-level questions
- **Professional Level:** Expert or graduate-level questions

## H 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>

1134	• <b>Phi-3.5-mini-instruct:</b>	<a href="https://huggingface.co/microsoft/Phi-3.5-mini-instruct">https://huggingface.co/microsoft/Phi-3.5-mini-instruct</a>
1135		
1136		
1137	• <b>Yi-1.5-6B-Chat:</b>	<a href="https://huggingface.co/01-ai/Yi-1.5-6B-Chat">https://huggingface.co/01-ai/Yi-1.5-6B-Chat</a>
1138		
1139	• <b>Qwen2.5-7B:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-7B">https://huggingface.co/Qwen/Qwen2.5-7B</a>
1140		
1141	• <b>vicuna-7b-v1.5:</b>	<a href="https://huggingface.co/lmsys/vicuna-7b-v1.5">https://huggingface.co/lmsys/vicuna-7b-v1.5</a>
1142		
1143	• <b>Qwen2-1.5B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2-1.5B-Instruct">https://huggingface.co/Qwen/Qwen2-1.5B-Instruct</a>
1144		
1145		
1146	• <b>phoenix-inst-chat-7b-v1.1:</b>	<a href="https://huggingface.co/FreedomIntelligence/phoenix-inst-chat-7b-v1.1">https://huggingface.co/FreedomIntelligence/phoenix-inst-chat-7b-v1.1</a>
1147		
1148		
1149	• <b>Qwen-Plus:</b>	<a href="https://huggingface.co/Qwen">https://huggingface.co/Qwen</a>
1150		
1151	• <b>GPT-3.5 turbo:</b>	<a href="https://openai.com/index/gpt-3-5-turbo-fine-tuning-and-api-updates/">https://openai.com/index/gpt-3-5-turbo-fine-tuning-and-api-updates/</a>
1152		
1153		
1154	• <b>doubao-1-5-pro-32k-250115:</b>	<a href="https://www.volcengine.com/product/doubao">https://www.volcengine.com/product/doubao</a>
1155		
1156	• <b>DeepSeek-V3-0324:</b>	<a href="https://huggingface.co/deepseek-ai/DeepSeek-V3-0324">https://huggingface.co/deepseek-ai/DeepSeek-V3-0324</a>
1157		
1158		
1159	• <b>doubao-1-5-lite-32k-250115:</b>	<a href="https://www.volcengine.com/product/doubao">https://www.volcengine.com/product/doubao</a>
1160		
1161	• <b>DeepSeek-R1-Distill-Llama-70B:</b>	<a href="https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-70B">https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-70B</a>
1162		
1163		
1164	• <b>Qwen2.5-72B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-72B-Instruct">https://huggingface.co/Qwen/Qwen2.5-72B-Instruct</a>
1165		
1166		
1167	• <b>Llama-3.1-70B-Instruct:</b>	<a href="https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct">https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct</a>
1168		
1169		
1170	• <b>DeepSeek-R1-Distill-Llama-8B:</b>	<a href="https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B">https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B</a>
1171		
1172		
1173	• <b>Athene-70B:</b>	<a href="https://huggingface.co/Nexusflow/Athene-70B">https://huggingface.co/Nexusflow/Athene-70B</a>
1174		
1175	• <b>gemma-2-27b-it:</b>	<a href="https://huggingface.co/google/gemma-2-27b-it">https://huggingface.co/google/gemma-2-27b-it</a>
1176		

## I 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  $\beta$  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:

$$\text{Calibrate\_Difficulty} = \text{Difficulty\_Label} \times 0.9^{t_i/\text{Reference\_Number}} \quad (8)$$

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.

## J 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

1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
~~1238~~

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  
Label:D

**MMLU-Pro:**

1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279

Example 1:  
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?  
A. Low-carbohydrate diet leading to nutritional deficiencies  
B. Brisk walk routine causing excessive physical exertion  
C. Medication side effects causing insulin fluctuations  
D. Herbal supplement causing hormonal imbalance  
E. Increased protein intake causing kidney strain  
F. Inadequate hydration from dietary changes  
G. Overconsumption of high-protein foods leading to weight gain  
H. Lack of fiber intake affecting metabolism  
I. Decrease in carbohydrate intake causing energy depletion  
J. Stress from lifestyle changes impacting health  
Label:C

1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296

Example 2:  
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.  
A. 0.1  
B. 0.2  
C. 0.3  
D. 0.4  
E. 0.5  
F. 0.6  
G. 0.7  
H. 0.8  
I. 0.9  
J. 1.0  
Label:J

1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
~~1320~~

**HellaSwag:**

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?  
A. 'Congratulations, Mark! I'm really proud of you! But I can't help feeling a bit overshadowed since I organized the event.'  
B. 'Wow, Mark! You totally deserve this! Yet, it's tough for me to celebrate when my efforts went unnoticed.'  
C. 'That was an amazing award, Mark! I'm happy for you! However, it stings that my contribution was overlooked.'  
D. 'I'm so thrilled for you, Mark! Your achievement is incredible! But it feels a little unfair that I didn't get a shoutout.'  
Label:C

1322  
  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
~~1360~~

1362  
1363  
1364  
  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
  
1403  
  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428

## K Examples of Model Generated Descriptions on Math Reasoning (Algebra)

### Task Description:

Evaluate the model's ability to solve algebra problems, including solving equations, simplifying expressions, and interpreting algebraic relationships.

### Query Description:

A clearly stated algebraic problem that includes sufficient details and conditions for solving. The question should test the model's ability to manipulate variables, apply algebraic principles, solve equations, or interpret graphs and systems of equations. The problem may involve concepts such as linear equations, quadratic equations, inequalities, exponents, or word problems requiring algebraic modeling.

### Label Format Description:

Includes one correct answer and several incorrect options. Incorrect options should reflect common mistakes, such as arithmetic errors, misunderstanding of algebraic rules, or misinterpretation of the problem. These options should be designed to test whether the model can reliably differentiate between valid and invalid approaches to solving the problem.

## L Prompt List

### LLM as Label Correctness Judge:

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

**Question:**  
{question}

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

Please note that the given question may be unsolvable, have a unique solution, multiple solutions, etc. Therefore, you should carefully analyze the correctness of the response to be checked based on the given question.

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

- 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.

- 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.
- 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.
- 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 response, as well as your confidence level in that judgment.

Your output should follow the template and example below:

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

**Example**  
Analyses:{Your detailed analyses}  
Judgement:1  
**Example End**

Now begin with "Analyses:"

### LLM as Comparison-based Label Correctness Judge:

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

**Question**  
{question}

1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498

1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548

```
###Candidate 1:
{{can1}}

###Candidate 2:
{{can2}}

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

You should begin your response with "
Correctness Analysis".
```

### LLM as Task Relevance Judge:

1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568

```
You are an expert who excels at
analyzing whether a given question
can be used to assess a specific
ability.

**Question:**
{{question}}

**Ability:**
{{ability}}

You should first carefully analyze what
abilities the given question can be
used to test.
Based on this analysis, compare it with
the given abilities.
After completing your analysis, please
```

```
provide your judgment on whether the
given question can be used to test
the given ability, as well as your
confidence in that judgment.

Your output should follow the template
below:
Analyses:{Your detailed analyses}
Judgement:
{output 0 if: You believe the given
question is completely unable to
test the given ability;
output 0.5 if: You believe the given
question is primarily meant to test
other abilities, but can also test
the given ability to some extent;
output 1 if: You believe the given
question primarily tests the given
ability.}

Now begin with "Analyses:"
```

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

```
You are a knowledgeable benchmark
creator.
Your task is to generate a creative
questions based on the provided Task
Description, Query Description,
Label Format Description, Generation
Guidelines, and Output Description
to help build a benchmark that
assesses the given task.

### Task Description:
{{task define}}

### Query Description:
{{query define}}

### Label Format Description:
{{label define}}

### Generation Guidelines:
1. Analyze the given task and think step
-by-step about the content needed to
construct the question, begin with
"Analyses:".
2. Generate the question content, begin
with "Question:".
3. Generate the right label strictly
following the Label Format
Description, begin with "Right Label
:".

### Output Description:
Strictly follow the template below to
```

1635 generate your sample.  
 1636 **\*\*Template\*\***  
 1637 **##Analyses:##** {{You analyze the provided  
 1638 attributes and outline the process  
 1639 for constructing the question to be  
 1640 generated.}}  
 1641 **##Question:##** {{Your generated question  
 1642 content}}  
 1643 **##Right Label:##**{{Your label to the  
 1644 question, strictly following the  
 1645 Label Format Description}}  
 1646 **\*\*Template End\*\***  
 1647  
 1648 Attention: You need to **\*\*strictly follow**  
 1649 the **template\*\*** and don't generate  
 1650 any other contents. Begin your  
 1651 response with "**##Analyses:##** "

Depth of Required Knowledge is Surface-  
 level 1702  
 1703  
 Depth of Required Knowledge is In-depth 1704  
 Depth of Required Knowledge is  
 Comprehensive 1705  
 1706

### Prompt of BENCHMARKER:

1707  
 1708  
 1709 You are a knowledgeable benchmark  
 1710 creator.  
 1711  
 1712 Your task is to generate a creative  
 1713 question based on the provided Task  
 1714 Description, Query Description,  
 1715 Label Format Description, General  
 1716 Attributes Descriptions, Difficulty  
 1717 Strategies Description, Generation  
 1718 Guidelines, and Output Description  
 1719 to help build a benchmark that  
 1720 assesses the given task.  
 1721

## M Examples of the Generated Difficulty Strategies

### MATH:

1652  
 1653  
 1654  
 1655  
 1656 Strategy 1:  
 1657 Complexity of Biological Concept is  
 1658 Basic  
 1659  
 1660 Complexity of Biological Concept is  
 1661 Intermediate  
 1662 Complexity of Biological Concept is  
 1663 Advanced  
 1664  
 1665 Strategy 2:  
 1666 Required Reasoning Steps set as Single-  
 1667 step  
 1668 Required Reasoning Steps set as Multi-  
 1669 step (2-3 steps)  
 1670 Required Reasoning Steps set as Multi-  
 1671 step (4-6 steps)  
 1672 Required Reasoning Steps set as More  
 1673 than 6 steps  
 1674  
 1675 Strategy 3:  
 1676 Familiarity with the Topic is Common  
 1677 Familiarity with the Topic is Uncommon  
 1678 Familiarity with the Topic is Rare  
 1679  
 1680 Strategy 4:  
 1681 Type of Biological Data Analysis is  
 1682 Qualitative  
 1683 Type of Biological Data Analysis is  
 1684 Quantitative  
 1685 Type of Biological Data Analysis is  
 1686 Advanced Data Interpretation  
 1687  
 1688 Strategy 5:  
 1689 Application of Concepts is Direct  
 1690 Application of Concepts is Modified  
 1691 Application of Concepts is Novel  
 1692  
 1693 Strategy 6:  
 1694 Integration Across Biological  
 1695 Disciplines is Single-discipline  
 1696 Integration Across Biological  
 1697 Disciplines is Cross-disciplinary  
 1698 Integration Across Biological  
 1699 Disciplines is Interdisciplinary  
 1700  
 1701 Strategy 7:

1722 **### Overall Task Description:**  
 1723 {{original task}}  
 1724  
 1725 **### Detailed Task Description:**  
 1726 {{task define}}  
 1727  
 1728 **### Query Description:**  
 1729 {{query define}}  
 1730  
 1731 **### Label Format Description:**  
 1732 {{label define}}  
 1733  
 1734 **### General Attributes Description:**  
 1735 You can refer to the following  
 1736 attributes and their corresponding  
 1737 values to construct questions, which  
 1738 means the questions you generate  
 1739 should ideally align with some of  
 1740 these attributes.  
 1741 Please note, if you find any conflicting  
 1742 or confusing parts among the  
 1743 attributes listed, you may disregard  
 1744 them.  
 1745  
 1746 {{attribute define}}  
 1747  
 1748 **### Difficulty Strategies Description:**  
 1749 Your generated questions should meet the  
 1750 following difficulty attribute  
 1751 requirements. If you find conflicts  
 1752 among these requirements, you may  
 1753 choose to selectively ignore them.  
 1754  
 1755 {{difficulty attribute define}}  
 1756  
 1757 **### Difficulty Description:**  
 1758 The following are some samples (0 or  
 1759 several).  
 1760 Please ensure that the difficulty level  
 1761 of the samples you generate is  
 1762 harder than these examples.  
 1763 The samples you generate should aim to  
 1764 assess different knowledge and  
 1765 skills compared to the given samples  
 1766 .  
 1767 The format of given samples are not what  
 1768  
 1769  
 1770  
 1771

1772 you should follow.

1773 **\*\*Please ensure that the sample you**

1774 **create differ substantially from the**

1775 **following samples, so as to**

1776 **maintain diversity in the resulting**

1777 **benchmark.\*\***

1778 **{{demonstrations}}**

1779

1780

1781 **### Generation Guidelines:**

1782 **\*\*Stage 1: Analyze\*\***

1783 **In this stage, you should analyze**

1784 **following the steps below and begin**

1785 **with "###Analyses:##". \*\*You need to**

1786 **clearly articulate the analysis**

1787 **content for each step\*\*, which means**

1788 **after completing Stage 1, you**

1789 **should have already produced a**

1790 **question that meets the requirements**

1791 **along with a correct and unique**

1792 **answer.**

1793 **1-1. Analyze the general attributes,**

1794 **difficulty attributes and difficulty**

1795 **description, and think step-by-step**

1796 **about the content needed to**

1797 **construct the question. \*\*Please use**

1798 **your imagination and avoid any**

1799 **obvious overlap with the given**

1800 **samples, either in the specific**

1801 **knowledge points being tested or in**

1802 **the format.\*\***

1803 **1-2. Start by drafting your question. If**

1804 **you discover any issues with the**

1805 **question or any overlapping parts**

1806 **between the generated question and**

1807 **the given samples during this**

1808 **process, feel free to revise it.**

1809 **1-3. Think through what the correct**

1810 **answer should be. If you discover**

1811 **any issues during this process,**

1812 **repeat the entire Stage 1 process**

1813 **from the beginning.**

1814 **1-4. Reevaluate your proposed question,**

1815 **answer to ensure that: the question**

1816 **meet the given attributes and**

1817 **Difficulty Description (you should**

1818 **compare the generated samples and**

1819 **given samples to verify this); the**

1820 **answer is both correct and unique.**

1821 **If it does not meet these criteria**

1822 **or you are not sure about this,**

1823 **repeat the entire Stage 1 process**

1824 **from the beginning.**

1825

1826 **\*\*Stage 2: Generate Sample\*\***

1827 **In this stage, you should give your**

1828 **generated sample in the right**

1829 **template based on the analyses above**

1830 **.**

1831 **2-1. Generate the question content,**

1832 **begin with "##Question:##".**

1833 **2-2. Generate a step-by-step reasoning**

1834 **process and the corresponding**

1835 **correct answer. Begin with "##**

1836 **Reasoning Path:##". If you find an**

1837 **issue with the question, return to**

1838 **Step 2-1 to regenerate the question.**

1839 **2-3. Generate the right label to the**

1840 **question strictly following the**

1841 **Label Format Description, begin with**

**"##Right Label:##".**

**### Output Description:**

**Strictly follow the template below to**

**generate your sample.**

**\*\*Template\*\***

**##Analyses:## {{You analyze the provided**

**attributes and outline the process**

**for constructing the question to be**

**generated.}}**

**##Question:## {{Your generated question**

**content}}**

**##Reasoning Path:## {{Your step-by-step**

**reasoning process}}**

**##Right Label:##{{Strictly follow the**

**Label Format Description to offer**

**the right label here}}**

**\*\*Template End\*\***

**Attention: You need to \*\*strictly follow**

**the template\*\* and don't generate**

**any other contents. Begin your**

**response with "##Analyses:##\n1-1. "**

1842

1843

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

1856

1857

1858

1859

1860

1861

1862

1863

1864

1865

**N Assessment Demands List**

1867

**Math Reasoning**

1868

**Subset Name: Prealgebra**

**Assessment Demands: Construct Prealgebra**

**test question with exactly one**

**correct answer for each. Use \boxed**

**{ } to denote the correct label.**

**Subset Name: Algebra**

**Assessment Demands: Construct Algebra**

**test question with exactly one**

**correct answer for each. Use \boxed**

**{ } to denote the correct label.**

**Subset Name: Number Theory**

**Assessment Demands: Construct Number**

**Theory test question with exactly**

**one correct answer for each. Use \**

**boxed{ } to denote the correct label.**

**Subset Name: Counting & Probability**

**Assessment Demands: Construct Counting &**

**Probability test question with**

**exactly one correct answer for each.**

**Use \boxed{ } to denote the correct**

**label.**

**Subset Name: Geometry**

**Assessment Demands: Construct Geometry**

**test question with exactly one**

**correct answer for each. Use \boxed**

**{ } to denote the correct label.**

**Subset Name: Intermediate Algebra**

**Assessment Demands: Construct**

**Intermediate Algebra test question**

**with exactly one correct answer for**

**each. Use \boxed{ } to denote the**

**correct label.**

**Subset Name: Precalculus**

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

1890

1891

1892

1893

1894

1895

1896

1897

1898

1899

1900

1901

1902

1903

1904

1905

1906

1907

1908

1909  
1910  
1911  
1912  
1913  
  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978

Assessment Demands: Construct  
Precalculus test question with  
exactly one correct answer for each.  
Use \boxed{} to denote the correct  
label.

### Language Understanding

Subset Name: psychology  
Assessment Demands:This benchmark is  
designed to assess **psychology**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: philosophy  
Assessment Demands:This benchmark is  
designed to assess **philosophy**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: health  
Assessment Demands:This benchmark is  
designed to assess **health**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: history  
Assessment Demands:This benchmark is  
designed to assess **history**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: business  
Assessment Demands:This benchmark is  
designed to assess **business**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: physics  
Assessment Demands:This benchmark is  
designed to assess **physics**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: engineering  
Assessment Demands:This benchmark is

designed to assess **engineering**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: chemistry  
Assessment Demands:This benchmark is  
designed to assess **chemistry**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: math  
Assessment Demands:This benchmark is  
designed to assess **math**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: computer science  
Assessment Demands:This benchmark is  
designed to assess **computer  
science** abilities while  
simultaneously evaluating knowledge  
understanding and complex reasoning  
skills, using **ten multiple-choice  
questions** as the evaluation format  
. Use \boxed{} to denote the correct  
label.

Subset Name: biology  
Assessment Demands:This benchmark is  
designed to assess **biology**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: economics  
Assessment Demands:This benchmark is  
designed to assess **economics**  
abilities while simultaneously  
evaluating knowledge understanding  
and complex reasoning skills, using  
**ten multiple-choice questions** as  
the evaluation format. Use \boxed{}  
to denote the correct label.

Subset Name: law  
Assessment Demands:This benchmark is  
designed to assess **law** abilities  
while simultaneously evaluating  
knowledge understanding and complex  
reasoning skills, using **ten  
multiple-choice questions** as the  
evaluation format. Use \boxed{} to  
denote the correct label.

1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
~~2046~~  
2048

### Commonsense Reasoning

2049  
2050  
2051  
2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
~~2061~~

Subset Name: NLI  
Assessment Demands: The task is to evaluate the model's commonsense natural language inference ability, using **four multiple-choice questions** as the evaluation format. Specifically, each question should present a concrete scenario, and the model should select the most likely event from the options based on a series of inferences. Use `\boxed{}` to denote the correct label.