

BENCHHUB: A UNIFIED BENCHMARK SUITE FOR HOLISTIC AND CUSTOMIZABLE LLM EVALUATION

Anonymous authors

Paper under double-blind review

ABSTRACT

As large language models (LLMs) continue to advance, the need for up-to-date and well-organized benchmarks becomes increasingly critical. However, many existing datasets are scattered, difficult to manage, and make it challenging to perform evaluations tailored to specific needs or domains, despite the growing importance of domain-specific models in areas such as math or code. In this paper, we introduce BENCHHUB, a dynamic benchmark repository that empowers researchers and developers to evaluate LLMs effectively, with a focus on Korean and English. BENCHHUB aggregates and automatically classifies benchmark datasets from diverse domains, integrating 839k questions across 54 benchmarks. It is designed to support continuous updates and scalable data management, enabling flexible and customizable evaluation tailored to various domains or use cases. Through extensive experiments with various LLM families, we demonstrate that model performance varies significantly across domain-specific subsets, emphasizing the importance of domain-aware benchmarking. Furthermore, we extend BENCHHUB into 10 languages spanning resource levels. We believe BenchHub can encourage better dataset reuse, more transparent model comparisons, and easier identification of underrepresented areas in existing benchmarks, offering a critical infrastructure for advancing LLM evaluation research.

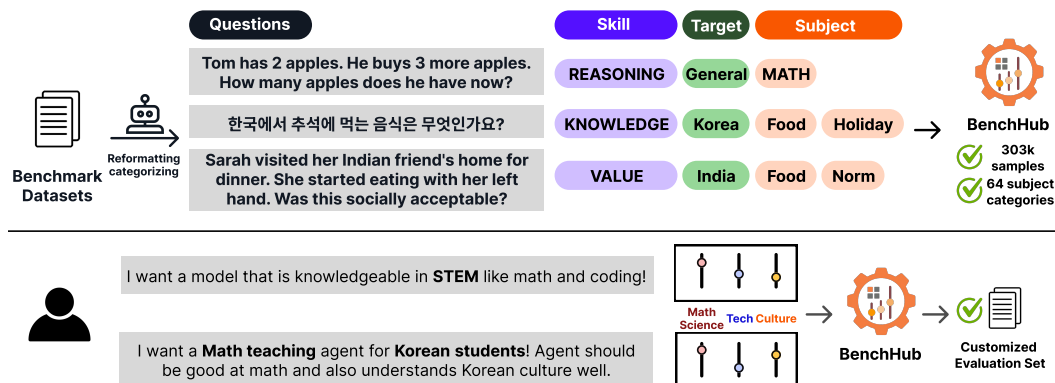


Figure 1: The concept of BENCHHUB. BENCHHUB automatically classifies and merges questions from existing benchmark datasets on a sample-wise basis. Through BENCHHUB, users can select test sets that align with their objectives and efficiently evaluate the models.

1 INTRODUCTION

Large language models (LLMs) have made remarkable strides, powering applications across diverse tasks, including research (Baek et al., 2025), industry (Chan et al., 2025), and everyday life (Chatterji et al., 2025). As their role varies with context—expanding into open-ended and high-stakes challenges ranging from utilizing external tools (Yang et al., 2023a) to making culturally sensitive decisions (Ki et al., 2025)—a new paradigm for LLM evaluation is essential. The key question is moving beyond

formulaic rankings toward a rigorous and comprehensive assessment of whether a model’s behavior aligns with the nuanced, custom objectives of specific users and applications.

In response, a wide range of evaluation efforts has emerged. On the one hand, holistic evaluation benchmarks (Liang et al., 2023; Ni et al., 2024) and leaderboards based on user preference (Chiang et al., 2024) or aggregated benchmarks (Aidar Myrzakhan, 2024) serve as popular community standards. While useful for broad comparisons, their aggregated scores obscure fine-grained strengths and weaknesses, often misaligning with the needs of specific applications Ribeiro et al. (2020). On the other hand, specialized benchmarks target narrow aspects, such as law (Li et al., 2024a), medical advice (Arora et al., 2025), and finance (Son et al., 2023), as well as specific tasks, including knowledge retrieval (Hendrycks et al., 2021a), reasoning (Cobbe et al., 2021; Zellers et al., 2019), and value alignment (Parrish et al., 2022; Ji et al., 2024). While these datasets capture critical capabilities, their vast, fragmented, and overlapping nature creates a chaotic landscape. For instance, in the mathematics domain, numerous benchmarks exist, such as MATH (Hendrycks et al., 2021b) and GSM8k (Cobbe et al., 2021), which in turn partially overlap with broader collections (e.g., MMLU (Hendrycks et al., 2021a)). This leaves researchers and practitioners with a dilemma: which benchmarks truly reflect their objective, and how can they compose a principled, customized evaluation suite tailored for diverse needs?

In this paper, we introduce BENCHHUB¹, a unified and customizable benchmark suite for holistic yet domain-aware LLM evaluation. BENCHHUB aggregates 839k questions from 54 benchmarks across 64 domains and 10 languages, mainly in English and Korean. We systematically categorize existing benchmarks by six dimensions: 1) tasks (e.g., mathematical reasoning), 2) answer formats (e.g., multiple-choice QA), 3) tool usage (i.e., language-only or requirements to external tools), 4) skills (i.e., knowledge, reasoning, or value/alignment), 5) coarse- and fine-grained subjects (e.g., STEM–mathematics), and 6) cultural-specificity (i.e., culturally specific or agnostic). This design facilitates users to dynamically construct their own evaluation sets tailored to their needs, moving beyond rigid, predefined test sets (Figure 1). To ensure long-term, dynamic scalability, we further train and release a categorization model that seamlessly integrates new, unseen benchmarks into BENCHHUB.

Using BENCHHUB, we evaluate 14 open LLMs and uncover a crucial insight: model rankings fluctuate substantially depending on benchmark compositions and domain focus. This finding highlights the central issue of benchmark composition bias, which can significantly distort interpretations of model performance. We further validate BENCHHUB through 5 real-world use cases—such as legal, educational, and cultural applications—showing how domain-aware evaluation alters conclusions about model superiority. We hope BENCHHUB provides a foundation for the community to move beyond monolithic leaderboards toward domain-aware, trustworthy, and customizable evaluation.

2 EXISTING LLM EVALUATION BENCHMARKS ARE SKEWED

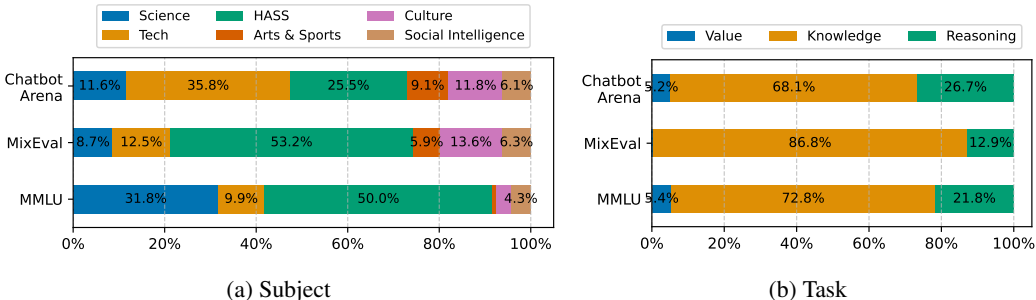


Figure 2: Data distribution of existing evaluation benchmarks.

¹We release our code and datasets at https://anonymous.4open.science/r/BenchHub_review-0A86; due to repository size limitations, only a subset is included, and the full dataset will be made available on Hugging Face following the anonymous review period.

108 What aspects do the commonly used multi-domain datasets evaluate, and how is the distribution of
 109 domains represented across these datasets? To answer this question, we classify three representative
 110 holistic benchmarks (*i.e.*, Chatbot Arena (Chiang et al., 2024), MixEval (Ni et al., 2024), and
 111 MMLU (Hendrycks et al., 2021a)) as multilabels using our fine-tuned classifiers (§ 3) in terms of
 112 coarse-grained subjects (Figure 2a) and tasks (Figure 2b). Among them, Chatbot Arena includes
 113 only 25.5% of Humanities and Social Science (HASS) questions, while both MixEval and MMLU
 114 comprise more than half of HASS questions. In addition, MixEval includes fewer than 0.30% of
 115 value alignment tasks and mostly focuses on measuring knowledge. Such disparities may lead to
 116 biased findings, where models that excel in certain domains may appear to perform better overall,
 117 potentially skewing the evaluation results.

118 Moreover, these biases are not limited to cross-
 119 benchmark comparisons but can also manifest
 120 within multilingual contexts. Figure 3 and Figure
 121 12 illustrate data distributions of MMLU
 122 series datasets in 5 languages classified by the
 123 model (§ 3) in terms of coarse-grained subjects.
 124 For instance, MMLU in English emphasizes
 125 HASS, whereas Korean MMLU (KMMLU) (Son
 126 et al., 2025b) comprises 76.1% of STEM (Sci-
 127 ence, Technology, Engineering, and Mathemat-
 128 ics) questions. This variation complicates the in-
 129 terpretation of performance differences, as it is
 130 challenging to discern whether the performance
 131 degradations in non-English are due to language
 132 proficiency or domain-specific knowledge.

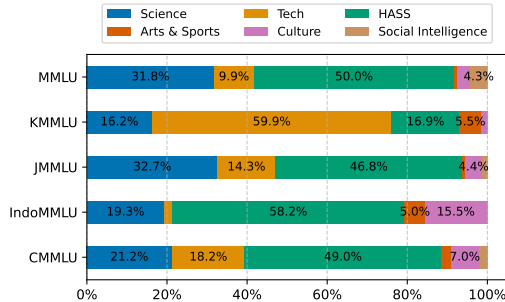


Figure 3: Data distribution of MMLU series in English, Korean, Japanese, Indonesian, and Chinese, respectively

132 Hence, instead of recklessly adopting existing holistic benchmarks, we recommend carefully selecting
 133 the benchmark suites for a reliable evaluation.

135 3 BENCHHUB

136 Consider a user who wants to determine “Which model excels at both mathematics and understanding
 137 culture?” As discussed in § 2, it remains unclear how to answer such specific, goal-oriented questions
 138 and how to construct their evaluation suite, as existing evaluation benchmarks (Hendrycks et al.,
 139 2021a; Liang et al., 2023; Ni et al., 2024) mainly provide general-purpose scores. To this end, we
 140 introduce BENCHHUB, a unified collection of LLM evaluation benchmarks across diverse domains.
 141 BENCHHUB integrates 54 benchmarks comprising 839k samples in 10 languages, with a primary
 142 focus on English and Korean as BENCHHUB-En and BENCHHUB-Ko, respectively. We design
 143 BENCHHUB around two core principles: 1) a fine-grained, multi-dimensional taxonomy to deconstruct
 144 model capabilities and 2) a fully automated pipeline to dynamically update and expand it with new
 145 datasets. In this section, we detail the taxonomy design (§ 3.1), the data curation (§ 3.2), the automated
 146 pipeline (§ 3.3), as well as interactive tools and utilities as a web-based platform (§ 3.4). Finally, we
 147 illustrate the multilingual extension of BENCHHUB—from English and Korean to eight additional
 148 languages—in § 3.5.

150 3.1 TAXONOMY

151 We annotate each dataset with six orthogonal dimensions: three dataset-level attributes—**task**, **answer**
 152 **format**, and **tool usage**— and three sample-level attributes—**skill**, **subject**, and **cultural-specificity**.
 153 The full scheme is illustrated in Appendix D.

154 Dataset-level attributes:

- 155 1. **Task** refers to the high-level family defined by the dataset authors (*e.g.*, mathematical reasoning,
 156 code generation, cultural understanding). This provides a general understanding of a dataset’s
 157 purpose. We assign it automatically from the dataset’s abstract or description using LLM inference.
- 158 2. **Answer format** specifies the expected response format: binary, multiple-choice QA (MCQA),
 159 short-form, free-form, open-ended (*e.g.*, story generation), and comparison (*e.g.*, determining
 160
 161

which response is better between A and B). This is crucial for selecting appropriate evaluation prompts and formats.

3. **Tool Usage** indicates whether a task requires language capabilities only (*language-only*) or interaction with external tools such as *e.g.*, code interpreters, web browsers, calculators (*requires externals tools*). This dimension supports agentic evaluation, where models must decide when and how to invoke external resources.

Sample-level attributes:

4. **Skill** captures the required ability to answer the question (*i.e.*, reasoning, knowledge, and value/alignment).
5. **Subject** denotes the knowledge domain. We define six coarse-grained categories—*Science, Technology, Humanities and Social Science (HASS), Arts & Sports, Culture, and Social Intelligence*—along with 64 sub-categories, by integrating various knowledge classification systems. Each sample may have multiple subject labels.
6. **Cultural-specificity** represents the cultural or geographical focus. Culturally agnostic items are labeled as *General*; otherwise, we assign a *Local* tag. This supports evaluation under culturally-aware evaluation (Singh et al., 2024).

3.2 DATASETS

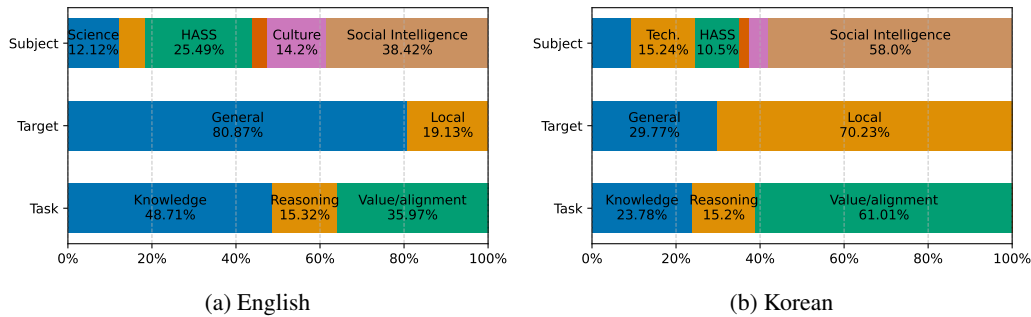


Figure 4: Data distribution of all datasets used in this paper by coarse-grained subjects, targets, and tasks. The English and Korean data include 250,940 and 144,331 questions each.

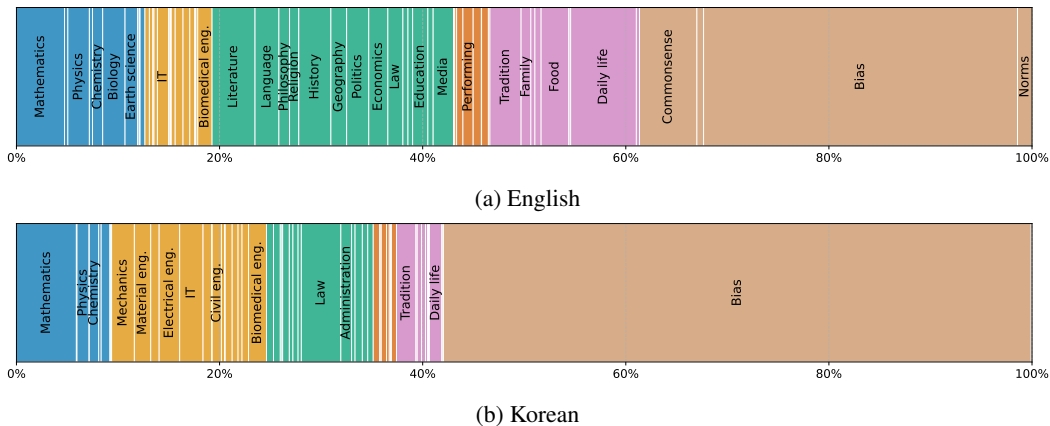


Figure 5: Fine-grained data distribution of all datasets used in this paper in terms of subjects

We apply this taxonomy to 54 benchmarks in 10 languages, mainly covering English and Korean and totaling over 839k instances. Figures 2 and 5 show the overall statistics of English and Korean datasets included in our benchmark. For English and Korean, we include 31 English and 12 Korean

language benchmarks with a total of 41 datasets.² We curate 1) general-purpose (*i.e.*, culturally agnostic) datasets commonly used by holistic evaluation benchmarks (Ye et al., 2024; Ni et al., 2024) and 2) culture-specific datasets. We select English datasets spanning multiple cultures drawn from a recent survey (Pawar et al., 2024), curating over 300 papers and datasets regarding LLM cultural awareness. For Korean, where public resources are fewer than in English, we include most datasets released after 2022. Table 2 in the Appendix provides a complete list of the datasets.

3.3 AUTOMATED AND DYNAMIC EXPANSION

With benchmark datasets emerging at a rapid pace, it is crucial to flexibly manage them for holistic evaluation. To dynamically adapt to newly emerging datasets, we automate the entire dataset merging process using an LLM agent, which includes reformatting the datasets into our benchmark format and classifying each sample into categories. The processing pipeline for a newly introduced dataset is outlined as follows:

1. **Reformatting:** We first automatically parse, reformat, and map a new dataset to the standardized BENCHHUB scheme using an LLM-guided rule-based approach. If the dataset does not adhere to our predefined schema, an LLM agent (*e.g.*, GPT-4o or Gemini) is employed to map keys to the correct format.
2. **Metadata assignment:** The LLM agent extracts the meta-task description and infers the task, answer format, and tool usage from the dataset documentation (*e.g.*, abstract).
3. **Sample-level Categorization:** We then assign sample-level attributes (*i.e.*, skill, subject, and **cultural-specificity**) using a fine-tuned Qwen-2.5-7B model (BenchHub-Cat-7B).³
4. **Merging:** The processed and annotated dataset is seamlessly merged into the main collections, thereby producing the next BENCHHUB release.

This automated pipeline allows BENCHHUB to continuously expand and provide more comprehensive evaluations as new datasets emerge. While we acknowledge the incompleteness of LLM-based expansion, we provide an empirical discussion of the reliability and robustness of this automated process in Appendix E.2.

Tables 1 show the accuracies of the categorizer model in Sample-level Categorization.

Table 1: Accuracy of fine-tuned categorizer on Qwen-2.5-7b

Sample-level Attribute	Accuracy
Subject	0.871
Skill	0.967
Cultural-specificity	0.986

3.4 INTERACTIVE PLATFORM AND UTILITIES

To proliferate our structured data into actionable insights for researchers and practitioners, we release an interactive web-based platform (Figure 10) and code utilities. The web demo allows users to filter out datasets by any category combinations, inspect statistics, download their customized subsets, and propose new datasets via pull requests. The code utilities offer two main features:

1. **Dataset Loader:** It filters the dataset to include only the categories selected by the user. It also allows the user to choose between returning the entire selected dataset or a filtered version with overlapping entries (including near-duplicates) removed, which is useful since multiple aggregated datasets may contain overlapping samples.
2. **Citation Report Generator:** For the customized dataset returned to the user, it produces a LaTeX table of datasets with their sources and licenses, includes dataset statistics such as the number of instances, and provides a comprehensive citation list (*e.g.*, BibTeX entries) to ensure proper credit to dataset authors.

²We count the multilingual datasets—BLEnD (Myung et al., 2024) and CaLMQA (Arora et al., 2024)—in both.

³The model link will be added after the anonymous review period. Details on the training and validation of BenchHub-Cat-7B are provided in Appendix E.1.

For better reproducibility, we adopt HRET (Lee et al., 2025)⁴, enabling direct evaluations on BENCHHUB. Design and implementation details of the platform and code utilities appear in Appendix B.

3.5 MULTILINGUAL EXTENSION OF BENCHHUB

While we focus on two languages (*i.e.*, Korean and English), we highlight that BENCHHUB is a language-agnostic, flexible framework that can be easily extended to other languages. To empirically guide this extension, we present BenchHub-Multi-Cat-7B⁵, a multilingual categorizer supporting 10 languages—English (En); 3 high-resource (Arabic (Ar), German (De), Dutch (Nl)); 3 mid-resource (Indonesian (Id), Korean (Ko), Ukrainian (Uk)); 3 low-resource (Swahili (Sw), Nepali (Ne), Kyrgyz (Ky)). Our multilingual categorizer achieves an average accuracy of 77.5% on fine-grained subject categorizations for unseen, out-of-domain data. Furthermore, we introduce BENCHHUB-multilingual, which extends our benchmark suite to a total of 10 languages consisting of 13 datasets and 444,402 samples. We hope BENCHHUB-multilingual to serve as a foundational step for reliable LLM evaluations in non-English languages. The details of the training procedure and the datasets for each language are provided in the Appendix F.

4 HOLISTIC AND CUSTOMIZABLE EVALUATION USING BENCHHUB

4.1 DOMAIN-AWARE EVALUATION USING BENCHHUB

In this section, we empirically show why the categories provided by BENCHHUB are important for LLM evaluation: performance varies substantially across categories (§ 4.1.1). Consequently, the dataset’s category distribution strongly influences model scores and leaderboard rankings (§ 4.1.2).

4.1.1 IMPACT OF SUBJECT CATEGORY ON MODEL RANKINGS

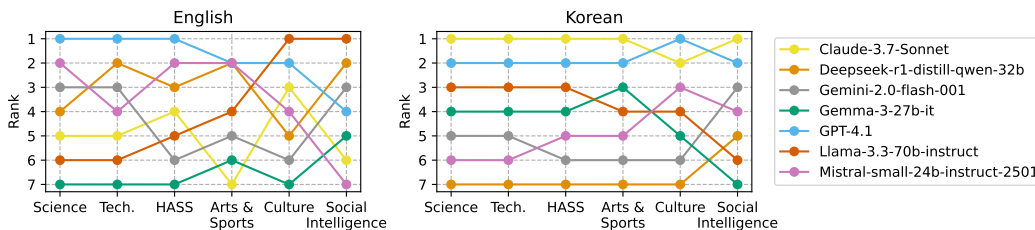


Figure 6: LLM evaluation ranking under BENCHHUB in terms of coarse-grained subjects

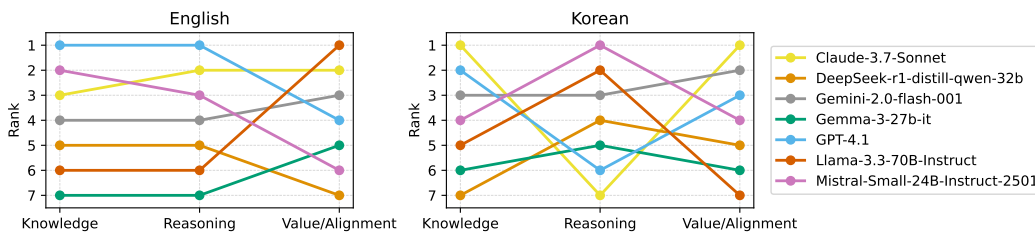


Figure 7: LLM evaluation ranking under BENCHHUB in terms of skills

In this section, we evaluate seven LLMs across diverse subjects using BENCHHUB. We select 6,644 and 6,485 examples for English and Korean, respectively. To manage the large number of fine-grained categories, we sample up to 150 examples per category, fully including categories with 100–150 samples and merging categories with fewer than 80 samples into a miscellaneous group within the same coarse-grained category. We extract the model’s intended answer from MCQA questions by

⁴HRET is an evaluation toolkit supporting multiple datasets, including BENCHHUB.

⁵The model link will be added after the anonymous review period.

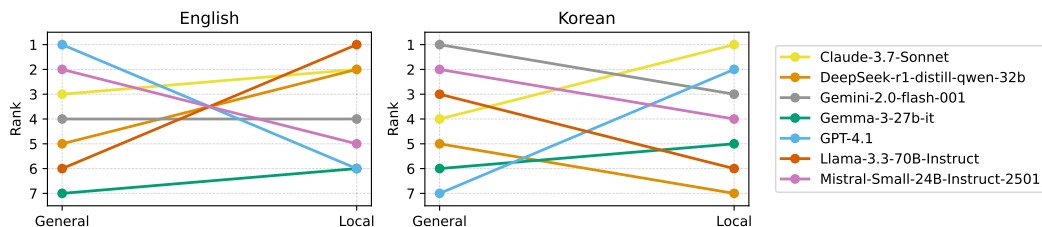


Figure 8: LLM evaluation ranking under BENCHHUB in terms of cultural-specificity

applying a set of regular expressions (Molfese et al., 2025), while using an LLM as a parser extractor for short-form questions⁶, similar to the approach in previous work (Ni et al., 2024).

We include one representative model from each commonly used LLM family. For proprietary models, we use GPT-4.1, Gemini-2.0-flash, and Claude 3.7 Sonnet⁷. Open models include Qwen-3-32b (Yang et al., 2025), DeepSeek-R1-Distill-Qwen-32B (DeepSeek-AI et al., 2025), Llama-3.3-70B (Grattafiori et al., 2024), Mistral-Small-24B-Instruct, and gemma-2-27b-it (Team et al., 2025).

Figures 6, 7, 8 present model rankings by subject, skill, and cultural-specificity categories, respectively. Our results show that the rankings fluctuate frequently, depending on the sample-level category. For example, Llama-3.3-70b ranks 6th in Science and Tech., but ranks as the top-performing model in Culture and Social Intelligence. This highlights the importance of domain-aware evaluation aligned with the evaluation context and objectives. The full results on the scores for each subject and model are in Table 17- 18 in the Appendix I.

4.1.2 IMPACT OF SAMPLING STRATEGIES ON MODEL RANKINGS

In this section, we empirically validate the influence of category distributions within evaluation benchmarks on model rankings. Since this requires experiments on large datasets for statistical validation, we include 14 open models ranging from 1B to 72B parameters. We test on a diverse set of English and Korean datasets, comprising 16,898 and 18,977 MCQA samples, respectively. The number of answer choices per MCQA sample varies between 3 and 18. We extract the model’s intended answer by applying a set of regular expressions (Molfese et al., 2025). The evaluated LLMs include:

- Qwen (Yang et al., 2024; 2025): Qwen2.5-72B-Instruct, Qwen3-1.7B, Qwen3-4B, Qwen3-8B, Qwen3-14B, Qwen3-32B
- DeepSeek (DeepSeek-AI et al., 2025): DeepSeek-R1-Distill-Qwen-14B, DeepSeek-R1-Distill-Qwen-32B
- Llama (Grattafiori et al., 2024): Llama-3.1-8B-Instruct, Llama-3.3-70B-Instruct
- Mistral: Mistral-Small-24B-Instruct-2501
- Gemma (Team et al., 2025): gemma-3-1b-it, gemma-3-4b-it, gemma-3-27b-it

To gauge the impact of data composition, we experiment under three sampling strategies with four setups, which are representatives of traditional approaches or emerging trends in LLM evaluations with a massive benchmark scale.

- **Random sampling:** Samples are drawn uniformly at random from the entire dataset collection, disregarding category proportions. Each sample has an equal chance of selection.
- **Stratified sampling:** Samples are drawn to ensure equal representation from each constituent dataset, preserving dataset-level balance rather than the overall distribution.
- **Sampling according to category distribution:** This strategy performs stratified sampling guided by fine-grained category distributions observed in existing holistic LLM benchmarks.

⁶We use GPT-4.1-nano as a parser extractor. Note that Ni et al. (2024) use GPT-3.5. The LLM parses and compares the extracted answer with the ground truth, without assessing answer quality.

⁷For GPT-4.1, we use GPT-4.1-2025-04-14 version. We directly call GPT-4.1 via the OpenAI API, while we use OpenRouter for Gemini-2.0-flash, and Claude 3.7 Sonnet.

We adopt the distributions derived from Chatbot Arena and MixEval, classified by our fine-tuned model (§ 3.3). The coarse-grained category distributions of these benchmarks are detailed in § 2.

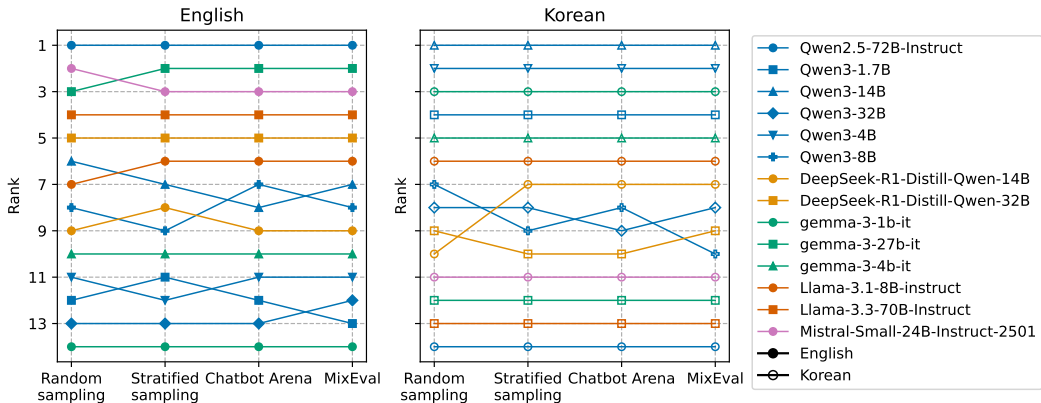


Figure 9: LLM ranking according to four sampling methods

We run 50 simulations per sampling setup, each selecting 5K questions. Model rankings within each setup follow normal distributions. Figure 9 visualizes LLM ranking changes across the four sampling setups. We use the Friedman test and the pairwise Wilcoxon test to statistically identify whether the sampling strategy affects the model ranking based on average accuracy. We observe a statistically significant difference across sampling strategies using the Friedman test ($p < 0.01$). Specifically, pairwise Wilcoxon signed-rank tests confirm that all pairs of sampling setups significantly differ in average, except for random sampling versus sampling according to MixEval distribution ($p < 0.01$). These findings underscore that category distribution and sampling strategy of data substantially affect LLM leaderboard rankings. We call on researchers and practitioners to carefully consider benchmark composition when evaluating LLMs. The full results for each subject and model are in Tables 19 and 20 in the Appendix I.

4.2 CUSTOMIZED EVALUATION USING BENCHHUB

In this section, we showcase how customized benchmark composition using BENCHHUB enables more targeted and meaningful evaluations tailored to real-world application scenarios. Here, we consider two use cases illustrated in Figure 1, and construct corresponding customized BENCHHUB as follows:

- (a) **STEM knowledge evaluation:** To identify the best-performing model with expertise in STEM domains, we select English datasets within BENCHHUB whose coarse-grained subjects are labeled as *Science* or *Technology*. To ensure balanced representation across individual datasets, the questions are drawn using a stratified sampling strategy at a dataset level.
- (b) **Math teaching agent for Korean students:** To evaluate Math teaching agents, we select Korean datasets comprising 1) math-related samples (*i.e.*, fine-grained categories are *Science/Math* or *Science/Statistics*), 2) education-related samples (*i.e.*, fine-grained category is *HASS/Education*), and 3) samples culturally specific to Korea (*i.e.*, **cultural-specificity** as ‘KR’). The final accuracy is computed as a weighted average of these subsets, with weights of 0.6, 0.1, and 0.3, respectively, reflecting their relative importance to the application.

Table 2 presents the rankings of LLMs under these customized benchmarks. We use the same set of models described in § 4.1.2. Notably, the model rankings differ substantially depending on the benchmark compositions, underscoring the practical need for tailored evaluations. The full results for each subject and model are in Table 21 in the Appendix I. We provide three additional real-world use cases (*i.e.*, legal chatbot, docent for Korean art, and counseling agent) and their corresponding model results in Appendix G.2.

Table 2: Top-5 LLMs evaluated by BENCHHUB in real-world application scenarios

Rank	(a) STEM knowledge evaluation (EN)		(b) Math teaching agent for Korean students (KO)	
	Customized	Stratified	Customized	Stratified
1	Qwen3-32B	gemma-3-1b-it	Qwen2.5-72B-Instruct	Qwen2.5-72B-Instruct
2	gemma-3-1b-it	Qwen3-32B	Mistral-Small-24B-Instruct-2501	Llama-3.3-70B-Instruct
3	Qwen3-1.7B	Qwen3-4B	gemma-3-27b-it	gemma-3-27b
4	Qwen3-4B	Qwen3-1.7B	Llama-3.3-70B-Instruct	Mistral-Small-24B-Instruct-2501
5	DeepSeek-R1-Distill-Qwen-14B	gemma-3-4b	DeepSeek-R1-Distill-Qwen-32B	DeepSeek-R1-Distill-Qwen-32B

5 RELATED WORK

As LLMs have become integral to real-world generative AI systems, the historical focus on benchmarks and leaderboards has matured into evaluation *science* (Weidinger et al., 2025). While LLM evaluation benchmarks primarily adopt a question-answering task as a default evaluation format, they have expanded their capabilities into diverse tasks, including long-form generation (Min et al., 2023), multilingual (Singh et al., 2024; Shafayat et al., 2024), multimodal (Fu et al., 2024), and complex reasoning tasks (Cobbe et al., 2021; Zellers et al., 2019), *inter alia*. This diversification reflects a growing recognition of the multifaceted capabilities and applications of LLMs.

Domain-specific Evaluation. Beyond general-purpose benchmarks, there has been a surge in domain-specific evaluation benchmarks targeting verticals such as healthcare and medicine (Hertzberg and Lokrantz, 2024; Matos et al., 2025; Rawat et al., 2024), law (Li et al., 2024a), science (Dinh et al., 2024), and financial (Zhang et al., 2025; Son et al., 2024a). These benchmarks enable more targeted assessment aligned with the unique requirements and challenges of each field. However, many domain-specific benchmarks lack the detail needed to compare specific skills or topics, and they often offer limited interoperability or consistency across benchmarks, making cross-benchmark comparison difficult. Complementing this trend, several large-scale benchmarks now aggregate tasks across multiple domains to facilitate robust, holistic evaluation of LLMs (Hendrycks et al., 2021a; Wang et al., 2024d; Taghanaki et al., 2024; Wang et al., 2022). However, it’s often unclear what the entire dataset actually evaluates, and thus lacks support for user-driven evaluation customization. In contrast, our paper proposes a framework that leverages existing benchmarks while enabling users to construct personalized, cross-domain evaluations tailored to their specific needs and contexts.

Dynamic Evaluation. Recent studies have identified inherent limitations of static datasets. Notably, issues such as data contamination, model overfitting to benchmarks, and insufficient human alignments have been highlighted (Yang et al., 2023b; Oren et al., 2024). This has spurred calls for a new discipline of *model metrology* focused on dynamic, adaptive, and robust evaluation frameworks (Saxon et al., 2024). Accordingly, dynamic and live evaluation is being conducted through various approaches: by synthetically generating evaluation data in real time (Zhang et al., 2024; Shashidhar et al., 2025); by incorporating human-in-the-loop platforms for periodic updates (Kiela et al., 2021; Chiang et al., 2024); or by regularly integrating new benchmark datasets (Ni et al., 2024; Jain et al., 2024). Our work extends this paradigm by offering a live benchmarking platform that automatically merges and recategorizes the benchmarks into a unified structure. This design makes our system more flexible and scalable for evaluating LLMs across diverse use cases.

Fine-grained Evaluation. Recent studies have shed light on the diversity of scenarios, contexts, and metrics in holistic evaluations. For example, (Wang et al., 2024a) critiqued over-reliance on single leaderboard rankings for evaluating AI fairness, advocating for multi-dimensional measurements. Similarly, (Liang et al., 2023) reformulated existing benchmarks into a format of diverse scenarios and adopted multiple metrics for a truly holistic assessment. Fine-grained evaluations, such as decomposing coarse scoring into skill-level scoring for alignment (Ye et al., 2024), facilitate richer and interpretable results. These advancements collectively underscore a paradigm shift from narrow, static benchmarks toward customizable, multi-faceted evaluations that better reflect the complex real-world capabilities and risks of LLMs. To support this shift, we propose a framework that enables question-level categorization across three core skills and 64 subject domains, offering a more fine-grained and interpretable evaluation.

To the best of our knowledge, BENCHHUB is the first to support domain-specific evaluation with fine-grained skill and subject categorization, while enabling dynamic updates through an automated integration pipeline for new benchmarks. We unify qualified benchmark datasets from diverse sources into a consistent structure and apply fine-grained categorization, enabling a holistic, interpretable evaluation pipeline that aligns closely with user-specific evaluation intents.

6 CONCLUSION

The rapid advancements in large language models (LLMs) have highlighted the need for robust and comprehensive evaluation frameworks capable of addressing the diverse and expanding range of their applications. While existing benchmarks have provided valuable insights into specific domains and capabilities, the fragmented nature of these datasets and the lack of alignment with task-specific objectives often limit their utility in real-world scenarios. Moreover, the varying distributions of subject types within benchmarks can significantly influence the interpretation of model performance, further emphasizing the need for systematic and customizable evaluation methodologies.

In this work, we introduced BENCHHUB, a unified benchmark suite designed to address these challenges. By categorizing 839k questions from 54 benchmarks in 10 languages across skills, subjects, and **cultural-specificity types**, BENCHHUB enables users to filter and create tailored test sets for domain-aware and task-specific evaluations. The integration of a categorization model based on Qwen-2.5-7b automates this process, ensuring scalability and adaptability to new datasets. Our experiments demonstrated that model performance rankings can vary significantly depending on subject categories and dataset distributions, underscoring the critical role of benchmark composition in fair and meaningful evaluations.

We hope this work promotes domain-aware evaluation and careful benchmark design. BENCHHUB serves as a practical tool to support these goals across diverse users.

- **For developers and practitioners**, BENCHHUB serves as a tool for accurately assessing model capabilities in targeted scenarios. They can identify each model’s strengths and weaknesses and select the ones best suited to their specific applications.
- **For benchmark and evaluation researchers**, we hope that the unified structure of BENCHHUB facilitates comprehensive statistical analysis of the coverage of existing benchmarks across subjects and skills, helping to identify underrepresented areas and motivating the construction of new datasets that address existing gaps in current evaluation practices.

ETHICS STATEMENT

We used ChatGPT, Cursor, and GitHub Copilot to refine the writing and assist with coding. BENCHHUB is provided for evaluation purposes only and must not be used for training models. Because BENCHHUB aggregates datasets from multiple sources, users must review and comply with the license terms of each dataset.

REPRODUCIBILITY STATEMENT

The prompts we used, as well as the model configurations and training methods are described in Appendix E.1. We release the code for BENCHHUB pipeline via https://anonymous.4open.science/r/BenchHub_review-0A86. We will release the trained models and the full dataset via the HuggingFace.

REFERENCES

Zhiqiang Shen Aidar Myrzakhan, Sondos Mahmoud Bsharat. Open-llm-leaderboard: From multi-choice to open-style questions for llms evaluation, benchmark, and arena. *arXiv preprint arXiv:2406.07545*, 2024. URL <https://arxiv.org/abs/2406.07545>.

- 540 Rahul K. Arora, Jason Wei, Hicks Rebecca Soskin, Preston Bowman, Joaquin Quiñero-
541 Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex
542 Beutel, and Johannes Heidecke. HealthBench: Evaluating large language models
543 towards improved human health, 2025. URL [https://cdn.openai.com/pdf/
544 bd7a39d5-9e9f-47b3-903c-8b847ca650c7/healthbench_paper.pdf](https://cdn.openai.com/pdf/bd7a39d5-9e9f-47b3-903c-8b847ca650c7/healthbench_paper.pdf).
- 545 Shane Arora, Marzena Karpinska, Hung-Ting Chen, Ipsita Bhattacharjee, Mohit Iyyer, and Eunsol
546 Choi. CaLMQA: Exploring culturally specific long-form question answering across 23 languages.
547 *arXiv preprint arXiv:2406.17761*, 2024. URL <https://arxiv.org/abs/2406.17761>.
- 549 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,
550 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language
551 models. *arXiv preprint arXiv:2108.07732*, 2021. URL [https://arxiv.org/abs/2108.
552 07732](https://arxiv.org/abs/2108.07732).
- 553 Axolotl AI. Axolotl: Scalable fine-tuning framework for llms. [https://axolotl-ai-cloud.
554 github.io/axolotl/](https://axolotl-ai-cloud.github.io/axolotl/), 2025. Github.
- 556 Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. ResearchAgent: It-
557 erative research idea generation over scientific literature with large language models. In Luis
558 Chiruzzo, Alan Ritter, and Lu Wang, editors, *Proceedings of the 2025 Conference of the Na-*
559 *tions of the Americas Chapter of the Association for Computational Linguistics: Human Lan-*
560 *guage Technologies (Volume 1: Long Papers)*, pages 6709–6738, Albuquerque, New Mexico,
561 April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL
562 <https://aclanthology.org/2025.naacl-long.342/>.
- 563 Yonatan Bisk, Rowan Zellers, Ronan Le bras, Jianfeng Gao, and Yejin Choi. PIQA: Reasoning
564 about physical commonsense in natural language. *Proceedings of the AAAI Conference on*
565 *Artificial Intelligence*, 34(05):7432–7439, Apr. 2020. doi: 10.1609/aaai.v34i05.6239. URL [https:
566 //ojs.aaai.org/index.php/AAAI/article/view/6239](https://ojs.aaai.org/index.php/AAAI/article/view/6239).
- 568 Alan Chan, Kevin Wei, Sihao Huang, Nitarshan Rajkumar, Elija Perrier, Seth Lazar, Gillian K
569 Hadfield, and Markus Anderljung. Infrastructure for AI agents. *Transactions on Machine Learn-*
570 *ing Research*, 2025. ISSN 2835-8856. URL [https://openreview.net/forum?id=
571 Ckh17xN2R2](https://openreview.net/forum?id=Ckh17xN2R2).
- 572 Aaron Chatterji, Thomas Cunningham, David J Deming, Zoe Hitzig, Christopher Ong, Carl Yan
573 Shan, and Kevin Wadman. How people use ChatGPT, September 2025. URL [http://www.
574 nber.org/papers/w34255](http://www.nber.org/papers/w34255).
- 576 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
577 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
578 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. URL [https://
579 arxiv.org/abs/2107.03374](https://arxiv.org/abs/2107.03374).
- 580 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng
581 Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena:
582 An open platform for evaluating LLMs by human preference. In Ruslan Salakhutdinov, Zico
583 Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp,
584 editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of
585 *Proceedings of Machine Learning Research*, pages 8359–8388. PMLR, 21–27 Jul 2024. URL
586 <https://proceedings.mlr.press/v235/chiang24b.html>.
- 587 Yu Ying Chiu, Liwei Jiang, Bill Yuchen Lin, Chan Young Park, Shuyue Stella Li, Sahithya Ravi,
588 Mehar Bhatia, Maria Antoniak, Yulia Tsvetkov, Vered Schwartz, et al. CulturalBench: a robust,
589 diverse and challenging benchmark on measuring the (lack of) cultural knowledge of LLMs. *arXiv*
590 *preprint arXiv:2410.02677*, 2024. URL <https://arxiv.org/abs/2410.02677>.
- 592 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
593 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
arXiv preprint arXiv:1803.05457, 2018. URL <https://arxiv.org/abs/1803.05457>.

- 594 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
595 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
596 math word problems. *arXiv preprint arXiv:2110.14168*, 2021. URL [https://arxiv.org/
597 abs/2110.14168](https://arxiv.org/abs/2110.14168).
- 598
599 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu,
600 Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,
601 Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao
602 Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
603 Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,
604 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding,
605 Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang
606 Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong,
607 Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao,
608 Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang,
609 Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang,
610 Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L.
611 Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang,
612 Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng
613 Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng
614 Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan
615 Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang,
616 Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen,
617 Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li,
618 Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang,
619 Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan,
620 Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia
621 He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong
622 Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha,
623 Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang,
624 Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li,
625 Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen
626 Zhang. DeepSeek-R1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv
627 preprint arXiv:2501.12948*, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 628 Tu Anh Dinh, Carlos Mullov, Leonard Bärmann, Zhaolin Li, Danni Liu, Simon Reiß, Jueun
629 Lee, Nathan Lerzer, Jianfeng Gao, Fabian Peller-Konrad, Tobias Röddiger, Alexander Waibel,
630 Tamim Asfour, Michael Beigl, Rainer Stiefelwagen, Carsten Dachsbacher, Klemens Böhm,
631 and Jan Niehues. SciEx: Benchmarking large language models on scientific exams with hu-
632 man expert grading and automatic grading. In Yaser Al-Onaizan, Mohit Bansal, and Yun-
633 Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural
634 Language Processing*, pages 11592–11610, Miami, Florida, USA, November 2024. Associa-
635 tion for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.647. URL <https://aclanthology.org/2024.emnlp-main.647/>.
- 636
637 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
638 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. MME: A comprehensive evaluation
639 benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2024. URL
640 <https://arxiv.org/abs/2306.13394>.
- 641
642 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad
643 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,
644 Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev,
645 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru,
646 Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak,
647 Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu,
648 Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle
649 Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego
650 Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,
651 Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel

648 Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon,
649 Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan
650 Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet,
651 Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,
652 Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie
653 Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua
654 Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak,
655 Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley
656 Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence
657 Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas
658 Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri,
659 Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie
660 Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes
661 Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne,
662 Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal
663 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
664 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
665 Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie
666 Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana
667 Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie,
668 Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon
669 Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan,
670 Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas
671 Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami,
672 Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti,
673 Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier
674 Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia,
675 Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen
676 Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe
677 Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya
678 Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei
679 Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu,
680 Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit
681 Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury,
682 Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer,
683 Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu,
684 Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido,
685 Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changan Wang, Changkyu
686 Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer,
687 Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu,
688 Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc
689 Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily
690 Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers,
691 Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank
692 Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee,
693 Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan,
694 Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph,
695 Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog,
696 Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James
697 Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny
698 Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings,
699 Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai
700 Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik
701 Veeraraghavan, Kelly Michelen, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle
Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng
Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish
Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim
Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle
Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang,

- 702 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam,
703 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier,
704 Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia
705 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro
706 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani,
707 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu
708 Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey,
709 Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak
710 Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan,
711 Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng
712 Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang
713 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
714 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng,
715 Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez,
716 Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim
717 Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez,
718 Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
719 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable,
720 Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun
721 Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu,
722 Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef
723 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models.
arXiv preprint arXiv:2407.21783, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 724 Zhicheng Guo, Sijie Cheng, Hao Wang, Shihao Liang, Yujia Qin, Peng Li, Zhiyuan Liu, Maosong
725 Sun, and Yang Liu. StableToolBench: Towards stable large-scale benchmarking on tool learning of
726 large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of
727 the Association for Computational Linguistics: ACL 2024*, pages 11143–11156, Bangkok, Thailand,
728 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.664.
729 URL <https://aclanthology.org/2024.findings-acl.664/>.
- 730 Serhii Hamotskyi, Anna-Izabella Levbarg, and Christian Häning. Eval-UA-tion 1.0: Benchmark for
731 evaluating Ukrainian (large) language models. In Mariana Romanyshyn, Nataliia Romanyshyn,
732 Andrii Hlybovets, and Oleksii Ignatenko, editors, *Proceedings of the Third Ukrainian Natural
733 Language Processing Workshop (UNLP) @ LREC-COLING 2024*, pages 109–119, Torino, Italia,
734 May 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.unlp-1.13/>.
- 735 Md Arid Hasan, Maram Hasanain, Fatema Ahmad, Sahinur Rahman Laskar, Sunaya Upadhyay,
736 Vrunda N Sukhadia, Mucahid Kutlu, Shammur Absar Chowdhury, and Firoj Alam. NativQA:
737 Multilingual culturally-aligned natural query for LLMs. *arXiv preprint arXiv:2407.09823*, 2024.
738 URL <https://arxiv.org/abs/2407.09823>.
- 740 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
741 Steinhardt. Measuring massive multitask language understanding. In *International Confer-
742 ence on Learning Representations*, 2021a. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.
- 744 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
745 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math
746 dataset. In J. Vanschoren and S. Yeung, editors, *Proceedings of the Neural Information
747 Processing Systems Track on Datasets and Benchmarks*, volume 1, 2021b. URL https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/be83ab3ecd0db773eb2dc1b0a17836a1-Paper-round2.pdf.
- 750 Niclas Hertzberg and Anna Lokrantz. MedQA-SWE - a clinical question & answer dataset for
751 Swedish. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani
752 Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference on
753 Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages
754 11178–11186, Torino, Italia, May 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.975/>.

- 756 Faris Hijazi, Somayah AlHarbi, Abdulaziz AlHussein, Harethah Abu Shairah, Reem AlZahrani,
757 Hebah AlShamlan, Omar Knio, and George Turkiyyah. Arablegaleval: A multitask benchmark
758 for assessing arabic legal knowledge in large language models. In *The Second Arabic Natural*
759 *Language Processing Conference*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=3EHYXqKKLA)
760 [3EHYXqKKLA](https://openreview.net/forum?id=3EHYXqKKLA).
- 761 Pin-Lun Hsu, Yun Dai, Vignesh Kothapalli, Qingquan Song, Shao Tang, Siyu Zhu, Steven Shimizu,
762 Shivam Sahni, Haowen Ning, and Yanning Chen. Liger kernel: Efficient triton kernels for llm
763 training. *arXiv preprint arXiv:2410.10989*, 2024. URL [https://arxiv.org/abs/2410.](https://arxiv.org/abs/2410.10989)
764 [10989](https://arxiv.org/abs/2410.10989).
- 766 Jing Huang, Zhengxuan Wu, Christopher Potts, Mor Geva, and Atticus Geiger. RAVEL: Evaluating
767 interpretability methods on disentangling language model representations. In Lun-Wei Ku, Andre
768 Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association*
769 *for Computational Linguistics (Volume 1: Long Papers)*, pages 8669–8687, Bangkok, Thailand,
770 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.470.
771 URL <https://aclanthology.org/2024.acl-long.470/>.
- 772 Jafar Isbarov, Arofat Akhundjanova, Mammad Hajili, Kavsar Huseynova, Dmitry Gaynullin, Anar
773 Rzayev, Osman Tursun, Ilshat Saetov, Rinat Kharisov, Saule Belginova, Ariana Kenbayeva, Amina
774 Alisheva, Aizirek Turdubaeva, Abdullatif Köksal, Samir Rustamov, and Duygu Ataman. TUMLU:
775 A Unified and Native Language Understanding Benchmark for Turkic Languages, 2025. URL
776 <https://arxiv.org/abs/2502.11020>.
- 778 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
779 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
780 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 781 Jianchao Ji, Yutong Chen, Mingyu Jin, Wujiang Xu, Wenyue Hua, and Yongfeng Zhang. MoralBench:
782 Moral evaluation of LLMs. *arXiv preprint arXiv:2406.04428*, 2024. URL [https://arxiv.](https://arxiv.org/abs/2406.04428)
783 [org/abs/2406.04428](https://arxiv.org/abs/2406.04428).
- 785 Jiho Jin, Jiseon Kim, Nayeon Lee, Haneul Yoo, Alice Oh, and Hwaran Lee. KoBBQ: Korean bias
786 benchmark for question answering. *Transactions of the Association for Computational Linguistics*,
787 12:507–524, 2024. doi: 10.1162/tacl_a_00661. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.tacl-1.28/)
788 [tacl-1.28/](https://aclanthology.org/2024.tacl-1.28/).
- 789 Dayeon Ki, Rachel Rudinger, Tianyi Zhou, and Marine Carpuat. Multiple LLM agents debate
790 for equitable cultural alignment. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and
791 Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for*
792 *Computational Linguistics (Volume 1: Long Papers)*, pages 24841–24877, Vienna, Austria, July
793 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/
794 2025.acl-long.1210. URL <https://aclanthology.org/2025.acl-long.1210/>.
- 796 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vid-
797 gen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel,
798 Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams.
799 Dynabench: Rethinking benchmarking in NLP. In Kristina Toutanova, Anna Rumshisky, Luke
800 Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty,
801 and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the North American Chapter of*
802 *the Association for Computational Linguistics: Human Language Technologies*, pages 4110–4124,
803 Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.
804 324. URL <https://aclanthology.org/2021.naacl-main.324/>.
- 805 Eunsu Kim, Juyoung Suk, Philhoon Oh, Haneul Yoo, James Thorne, and Alice Oh. CLICk: A
806 benchmark dataset of cultural and linguistic intelligence in Korean. In Nicoletta Calzolari, Min-Yen
807 Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue, editors, *Proceedings*
808 *of the 2024 Joint International Conference on Computational Linguistics, Language Resources*
809 *and Evaluation (LREC-COLING 2024)*, pages 3335–3346, Torino, Italia, May 2024a. ELRA and
ICCL. URL <https://aclanthology.org/2024.lrec-main.296/>.

- 810 Yeeun Kim, Youngrok Choi, Eunkyung Choi, JinHwan Choi, Hai Jin Park, and Wonseok Hwang.
811 Developing a pragmatic benchmark for assessing Korean legal language understanding in
812 large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, edi-
813 tors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 5573–
814 5595, Miami, Florida, USA, November 2024b. Association for Computational Linguistics.
815 doi: 10.18653/v1/2024.findings-emnlp.319. URL [https://aclanthology.org/2024.
816 findings-emnlp.319/](https://aclanthology.org/2024.findings-emnlp.319/).
- 817 Hyunwoo Ko, Guijin Son, and Dasol Choi. Understand, solve and translate: Bridging the multilingual
818 mathematical reasoning gap. *arXiv preprint arXiv:2501.02448*, 2025. URL [https://arxiv.
819 org/abs/2501.02448](https://arxiv.org/abs/2501.02448).
- 820 Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis,
821 and Edward Grefenstette. The NarrativeQA reading comprehension challenge. *Transactions of the
822 Association for Computational Linguistics*, 6:317–328, 2018. doi: 10.1162/tacl_a_00023. URL
823 <https://aclanthology.org/Q18-1023/>.
- 824 Fajri Koto, Haonan Li, Sara Shatnawi, Jad Doughman, Abdelrahman Sadallah, Aisha Alraeesi,
825 Khalid Almubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, Nizar Habash, Preslav Nakov,
826 and Timothy Baldwin. ArabicMMLU: Assessing massive multitask language understanding in
827 Arabic. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association
828 for Computational Linguistics: ACL 2024*, pages 5622–5640, Bangkok, Thailand, August 2024.
829 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.334. URL [https://aclanthology.org/2024.
830 findings-acl.334/](https://aclanthology.org/2024.findings-acl.334/).
- 831 Sunjun Kweon, Byungjin Choi, Gyouk Chu, Junyeong Song, Daeun Hyeon, Sujin Gan, Jueon
832 Kim, Minkyu Kim, Rae Woong Park, and Edward Choi. KorMedMCQA: multi-choice question
833 answering benchmark for korean healthcare professional licensing examinations. *arXiv preprint
834 arXiv:2403.01469*, 2024. URL <https://arxiv.org/abs/2403.01469>.
- 835 Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
836 Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
837 Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
838 Petrov. Natural questions: A benchmark for question answering research. *Transactions of the
839 Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl_a_00276. URL
840 <https://aclanthology.org/Q19-1026/>.
- 841 Hanwool Lee, Dasol Choi, Sooyong Kim, Ilgyun Jung, Sangwon Baek, Guijin Son, Inseon Hwang,
842 Naeun Lee, and Seunghyeok Hong. Redefining evaluation standards: A unified framework for
843 evaluating the korean capabilities of language models, 2025. URL [https://arxiv.org/
844 abs/2503.22968](https://arxiv.org/abs/2503.22968).
- 845 Hwaran Lee, Seokhee Hong, Joonsuk Park, Takyoun Kim, Gunhee Kim, and Jung-woo Ha. KoSBI: A
846 dataset for mitigating social bias risks towards safer large language model applications. In Sunayana
847 Sitaram, Beata Beigman Klebanov, and Jason D Williams, editors, *Proceedings of the 61st Annual
848 Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 208–
849 224, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/
850 2023.acl-industry.21. URL <https://aclanthology.org/2023.acl-industry.21/>.
- 851 Jiyoung Lee, Minwoo Kim, Seungho Kim, Junghwan Kim, Seunghyun Won, Hwaran Lee, and Edward
852 Choi. KorNAT: LLM alignment benchmark for Korean social values and common knowledge.
853 In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for
854 Computational Linguistics: ACL 2024*, pages 11177–11213, Bangkok, Thailand, August 2024.
855 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.666. URL [https://aclanthology.org/2024.
856 findings-acl.666/](https://aclanthology.org/2024.findings-acl.666/).
- 857 Haitao Li, Junjie Chen, Jingli Yang, Qingyao Ai, Wei Jia, Youfeng Liu, Kai Lin, Yueyue Wu, Guozhi
858 Yuan, Yiran Hu, et al. LegalAgentBench: Evaluating LLM agents in legal domain. *arXiv preprint
859 arXiv:2412.17259*, 2024a. URL <https://arxiv.org/abs/2412.17259>.
- 860 Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion
861 Stoica. From live data to high-quality benchmarks: The arena-hard pipeline, April 2024b. URL
862 <https://lmsys.org/blog/2024-04-19-arena-hard/>.
- 863

- 864 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian
865 Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby
866 Yan, Ce Zhang, Christian Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas,
867 Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu
868 Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun,
869 Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang,
870 Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto,
871 Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui
872 Zhang, and Yuta Koreeda. Holistic evaluation of language models. *Transactions on Machine*
873 *Learning Research*, 2023. ISSN 2835-8856. URL [https://openreview.net/forum?](https://openreview.net/forum?id=iO4LZibEqW)
874 [id=iO4LZibEqW](https://openreview.net/forum?id=iO4LZibEqW). Featured Certification, Expert Certification, Outstanding Certification.
- 875 Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human
876 falsehoods. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of*
877 *the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
878 pages 3214–3252, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi:
879 10.18653/v1/2022.acl-long.229. URL [https://aclanthology.org/2022.acl-long.](https://aclanthology.org/2022.acl-long.229/)
880 [229/](https://aclanthology.org/2022.acl-long.229/).
- 881 Holy Lovenia, Rahmad Mahendra, Salsabil Maulana Akbar, Lester James Validad Miranda, Jennifer
882 Santoso, Elyanah Aco, Akhdan Fadhilah, Jonibek Mansurov, Joseph Marvin Imperial, Onno P.
883 Kampman, Joel Ruben Antony Moniz, Muhammad Ravi Shulthan Habibi, Frederikus Hudi,
884 Railey Montalan, Ryan Ignatius Hadiwijaya, Joanito Agili Lopo, William Nixon, Börje F. Karls-
885 son, James Jaya, Ryandito Diandaru, Yuze Gao, Patrick Amadeus Irawan, Bin Wang, Jan Chris-
886 tian Blaise Cruz, Chenxi Whitehouse, Ivan Halim Parmonangan, Maria Khelli, Wenyu Zhang,
887 Lucky Susanto, Reynard Adha Ryanda, Sonny Lazuardi Hermawan, Dan John Velasco, Muhammad
888 Dehan Al Kautsar, Willy Fitra Hendria, Yasmin Moslem, Noah Flynn, Muhammad Farid Adil-
889 lazuarda, Haochen Li, Johannes Lee, R. Damanhuri, Shuo Sun, Muhammad Reza Qorib, Amirbek
890 Djanibekov, Wei Qi Leong, Quyet V. Do, Niklas Muennighoff, Tanrada Pansuwan, Ilham Fir-
891 dausi Putra, Yan Xu, Tai Ngee Chia, Ayu Purwarianti, Sebastian Ruder, William Chandra Tjhi,
892 Peerat Limkonchotiwat, Alham Fikri Aji, Sedrick Keh, Genta Indra Winata, Ruochen Zhang,
893 Fajri Koto, Zheng Xin Yong, and Samuel Cahyawijaya. SEACrowd: A multilingual multi-
894 modal data hub and benchmark suite for Southeast Asian languages. In Yaser Al-Onaizan,
895 Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical*
896 *Methods in Natural Language Processing*, pages 5155–5203, Miami, Florida, USA, November
897 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.296. URL
898 <https://aclanthology.org/2024.emnlp-main.296/>.
- 899 João Matos, Shan Chen, Siena Kathleen V. Placino, Yingya Li, Juan Carlos Climent Pardo, Daphna
900 Idan, Takeshi Tohyama, David Restrepo, Luis Filipe Nakayama, José María Millet Pascual-Leone,
901 Guergana K Savova, Hugo Aerts, Leo Anthony Celi, An-Kwok Ian Wong, Danielle Bitterman,
902 and Jack Gallifant. WorldMedQA-V: a multilingual, multimodal medical examination dataset for
903 multimodal language models evaluation. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors,
904 *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 7203–7216,
905 Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-
906 89176-195-7. URL <https://aclanthology.org/2025.findings-naacl.402/>.
- 907 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
908 electricity? a new dataset for open book question answering. In Ellen Riloff, David Chiang,
909 Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empir-
910 ical Methods in Natural Language Processing*, pages 2381–2391, Brussels, Belgium, October-
911 November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260. URL
912 <https://aclanthology.org/D18-1260/>.
- 913 Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke
914 Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual preci-
915 sion in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceed-*
916 *ings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–
917 12100, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/
2023.emnlp-main.741. URL <https://aclanthology.org/2023.emnlp-main.741/>.

- 918 Francesco Maria Molfese, Luca Moroni, Luca Gioffrè, Alessandro Scirè, Simone Conia, and Roberto
919 Navigli. Right answer, wrong score: Uncovering the inconsistencies of llm evaluation in multiple-
920 choice question answering. *arXiv preprint arXiv:2503.14996*, 2025. URL <https://arxiv.org/abs/2503.14996>.
921
- 922 Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia
923 Shi, Evan Pete Walsh, Oyvind Taffjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia,
924 Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela,
925 Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. OLMoe:
926 Open mixture-of-experts language models. In *The Thirteenth International Conference on Learning
927 Representations*, 2025. URL <https://openreview.net/forum?id=xXTkbTBmqq>.
928
- 929 Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin, Rifki Afina Putri, Dimosthenis Antypas, Hsuvas
930 Borkakoty, Eunsu Kim, Carla Perez-Almendros, Abinew Ali Ayele, Víctor Gutiérrez-Basulto,
931 Yazmín Ibáñez García, Hwaran Lee, Shamsuddeen Hassan Muhammad, Kiwoong Park,
932 Anar Sabuhi Rzayev, Nina White, Seid Muhie Yimam, Mohammad Taher Pilehvar, Ned-
933 jma Ousidhoum, Jose Camacho-Collados, and Alice Oh. BLEND: A benchmark for llms
934 on everyday knowledge in diverse cultures and languages. In A. Globerson, L. Mackey,
935 D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural
936 Information Processing Systems*, volume 37, pages 78104–78146. Curran Associates, Inc., 2024.
937 URL [https://proceedings.neurips.cc/paper_files/paper/2024/file/
938 8eb88844dafefa92a26aaec9f3acad93-Paper-Datasets_and_Benchmarks_
939 Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/8eb88844dafefa92a26aaec9f3acad93-Paper-Datasets_and_Benchmarks_Track.pdf).
- 940 Tuan-Phong Nguyen, Simon Razniewski, Aparna Varde, and Gerhard Weikum. Extracting cul-
941 tural commonsense knowledge at scale. In *Proceedings of the ACM Web Conference 2023*,
942 WWW '23, page 1907–1917, New York, NY, USA, 2023. Association for Computing Machinery.
943 ISBN 9781450394161. doi: 10.1145/3543507.3583535. URL [https://doi.org/10.1145/
944 3543507.3583535](https://doi.org/10.1145/3543507.3583535).
- 945 Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig, and
946 Yang You. MixEval: Deriving wisdom of the crowd from LLM benchmark mixtures. In
947 A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors,
948 *Advances in Neural Information Processing Systems*, volume 37, pages 98180–98212. Curran Asso-
949 ciates, Inc., 2024. URL [https://proceedings.neurips.cc/paper_files/paper/
950 2024/file/b1f34d7b4a03a3d80be8e72eb430dd81-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/b1f34d7b4a03a3d80be8e72eb430dd81-Paper-Conference.pdf).
- 951 Jinu Nyachhyan, Mridul Sharma, Prajwal Thapa, and Bal Krishna Bal. Consolidating and developing
952 benchmarking datasets for the nepali natural language understanding tasks, 2025. URL [https:
953 //arxiv.org/abs/2411.19244](https://arxiv.org/abs/2411.19244).
- 954 Yonatan Oren, Nicole Meister, Niladri S. Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. Pro-
955 ving test set contamination in black-box language models. In *The Twelfth International Confer-
956 ence on Learning Representations*, 2024. URL [https://openreview.net/forum?id=
957 KS8mIvetg2](https://openreview.net/forum?id=KS8mIvetg2).
- 958 Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson,
959 Phu Mon Htut, and Samuel Bowman. BBQ: A hand-built bias benchmark for question answer-
960 ing. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Findings of the
961 Association for Computational Linguistics: ACL 2022*, pages 2086–2105, Dublin, Ireland, May
962 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.165. URL
963 <https://aclanthology.org/2022.findings-acl.165/>.
964
- 965 Sachin Pawar, Nitin Ramrakhiani, Anubhav Sinha, Manoj Apte, and Girish Palshikar. Why generate
966 when you can discriminate? a novel technique for text classification using language models. In
967 Yvette Graham and Matthew Purver, editors, *Findings of the Association for Computational Linguis-
968 tics: EAACL 2024*, pages 1099–1114, St. Julian’s, Malta, March 2024. Association for Computational
969 Linguistics. URL <https://aclanthology.org/2024.findings-eaACL.74/>.
- 970 Jan Pfister and Andreas Hotho. SuperGLEBer: German language understanding evaluation benchmark.
971 In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of
the North American Chapter of the Association for Computational Linguistics: Human Language*

- 972 *Technologies (Volume 1: Long Papers)*, pages 7904–7923, Mexico City, Mexico, June 2024.
973 Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.438. URL <https://aclanthology.org/2024.naacl-long.438/>.
974
- 975 Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. XCOPA: A multilingual dataset for causal commonsense reasoning. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.185. URL <https://aclanthology.org/2020.emnlp-main.185/>.
976
977
978
979
980
981
- 982 Rifki Afina Putri and Alice Oh. IDK-MRC: Unanswerable questions for Indonesian machine reading
983 comprehension. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of
984 the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6918–6933,
985 Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.465. URL [https://aclanthology.org/2022.
986 emnlp-main.465/](https://aclanthology.org/2022.emnlp-main.465/).
987
- 988 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru
989 Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li,
990 Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16,000+
991 real-world apis. *arXiv preprint arXiv:2307.16789*, 2023. URL [https://arxiv.org/abs/
992 2307.16789](https://arxiv.org/abs/2307.16789).
- 993 Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations
994 toward training trillion parameter models. In *SC20: International Conference for High Performance
995 Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020.
- 996 Abhinav Sukumar Rao, Akhila Yerukola, Vishwa Shah, Katharina Reinecke, and Maarten Sap.
997 NormAd: A framework for measuring the cultural adaptability of large language models. In Luis
998 Chiruzzo, Alan Ritter, and Lu Wang, editors, *Proceedings of the 2025 Conference of the Nations
999 of the Americas Chapter of the Association for Computational Linguistics: Human Language
1000 Technologies (Volume 1: Long Papers)*, pages 2373–2403, Albuquerque, New Mexico, April
1001 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. URL <https://aclanthology.org/2025.naacl-long.120/>.
1002
- 1003 Rajat Rawat, Hudson McBride, Rajarshi Ghosh, Dhiyaan Nirmal, Jong Moon, Dhruv Alamuri,
1004 Sean O’Brien, and Kevin Zhu. DiversityMedQA: A benchmark for assessing demographic
1005 biases in medical diagnosis using large language models. In Daryna Dementieva, Oana Ignat,
1006 Zhijing Jin, Rada Mihalcea, Giorgio Piatti, Joel Tetreault, Steven Wilson, and Jieyu Zhao, editors,
1007 *Proceedings of the Third Workshop on NLP for Positive Impact*, pages 334–348, Miami, Florida,
1008 USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.nlp4pi-1.
1009 29. URL <https://aclanthology.org/2024.nlp4pi-1.29/>.
- 1010 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
1011 Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof Q&A
1012 benchmark. In *First Conference on Language Modeling, 2024*. URL [https://openreview.
1013 net/forum?id=Ti67584b98](https://openreview.net/forum?id=Ti67584b98).
- 1014 Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. Beyond accuracy: Behavioral
1015 testing of NLP models with CheckList. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel
1016 Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational
1017 Linguistics*, pages 4902–4912, Online, July 2020. Association for Computational Linguistics. doi:
1018 10.18653/v1/2020.acl-main.442. URL [https://aclanthology.org/2020.acl-main.
1019 442/](https://aclanthology.org/2020.acl-main.442/).
- 1020 Muhammad Rizqullah, Ayu Purwarianti, and Alham Fikri Aji. QASiNa: Religious domain question
1021 answering using sirah nabawiyah. In *preprint / arXiv*, 2023. URL [https://arxiv.org/
1022 abs/2310.08102](https://arxiv.org/abs/2310.08102).
1023
- 1024 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: an adversarial
1025 winograd schema challenge at scale. *Commun. ACM*, 64(9):99–106, August 2021. ISSN 0001-0782. doi: 10.1145/3474381. URL <https://doi.org/10.1145/3474381>.

- 1026 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Common-
1027 sense reasoning about social interactions. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan,
1028 editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*
1029 *and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*,
1030 pages 4463–4473, Hong Kong, China, November 2019. Association for Computational Linguistics.
1031 doi: 10.18653/v1/D19-1454. URL <https://aclanthology.org/D19-1454/>.
- 1032 Michael Saxon, Ari Holtzman, Peter West, William Yang Wang, and Naomi Saphra. Benchmarks as
1033 microscopes: A call for model metrology. In *First Conference on Language Modeling*, 2024. URL
1034 <https://openreview.net/forum?id=bttKwCZDkm>.
- 1035 Sheikh Shafayat, Eunsu Kim, Juhyun Oh, and Alice Oh. Multi-fact: Assessing factuality of multilin-
1036 gual llms using factscore, 2024. URL <https://arxiv.org/abs/2402.18045>.
- 1037
1038 Sumuk Shashidhar, Clémentine Fourier, Alina Lozovskia, Thomas Wolf, Gokhan Tur, and
1039 Dilek Hakkani-Tür. Yourbench: Easy custom evaluation sets for everyone. *arXiv preprint*
1040 *arXiv:2504.01833*, 2025. URL <https://arxiv.org/abs/2504.01833>.
- 1041
1042 Weiyang Shi, Ryan Li, Yutong Zhang, Caleb Ziems, Sunny Yu, Raya Horesh, Rogério Abreu De
1043 Paula, and Diyi Yang. CultureBank: An online community-driven knowledge base towards
1044 culturally aware language technologies. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung
1045 Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages
1046 4996–5025, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
1047 doi: 10.18653/v1/2024.findings-emnlp.288. URL <https://aclanthology.org/2024.findings-emnlp.288/>.
- 1048
1049 Shivalika Singh, Angelika Romanou, Clémentine Fourier, David I. Adelani, Jian Gang Ngui, Daniel
1050 Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Raymond
1051 Ng, Shayne Longpre, Wei-Yin Ko, Sebastian Ruder, Madeline Smith, Antoine Bosselut, Alice
1052 Oh, Andre F. T. Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee,
1053 Beyza Ermis, and Sara Hooker. Global MMLU: Understanding and addressing cultural and
1054 linguistic biases in multilingual evaluation. *arXiv preprint arXiv:2412.03304*, 2024. URL <https://arxiv.org/abs/2412.03304>.
- 1055
1056 Guijin Son, Hanearl Jung, Moonjeong Hahm, Keonju Na, and Sol Jin. Beyond classification:
1057 Financial reasoning in state-of-the-art language models. *arXiv preprint arXiv:2305.01505*, 2023.
1058 URL <https://arxiv.org/abs/2305.01505>.
- 1059
1060 Guijin Son, Hyunjun Jeon, Chami Hwang, and Hanearl Jung. KRX bench: Automating finan-
1061 cial benchmark creation via large language models. In Chung-Chi Chen, Xiaomo Liu, Udo
1062 Hahn, Armineh Nourbakhsh, Zhiqiang Ma, Charese Smiley, Veronique Hoste, Sanjiv Ran-
1063 jan Das, Manling Li, Mohammad Ghassemi, Hen-Hsen Huang, Hiroya Takamura, and Hsin-
1064 Hsi Chen, editors, *Proceedings of the Joint Workshop of the 7th Financial Technology and*
1065 *Natural Language Processing, the 5th Knowledge Discovery from Unstructured Data in Fi-*
1066 *ncial Services, and the 4th Workshop on Economics and Natural Language Processing*,
1067 pages 10–20, Torino, Italia, May 2024a. Association for Computational Linguistics. URL
<https://aclanthology.org/2024.finnlp-1.2/>.
- 1068
1069 Guijin Son, Hanwool Lee, Suwan Kim, Huiseo Kim, Jae cheol Lee, Je Won Yeom, Jihyu Jung,
1070 Jung woo Kim, and Songseong Kim. HAE-RAE bench: Evaluation of Korean knowledge in
1071 language models. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci,
1072 Sakriani Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference*
1073 *on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages
1074 7993–8007, Torino, Italia, May 2024b. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.704/>.
- 1075
1076 Guijin Son, Hyunwoo Ko, and Dasol Choi. Multi-step reasoning in Korean and the emergent
1077 mirage. In Vinodkumar Prabhakaran, Sunipa Dev, Luciana Benotti, Daniel Hershcovich, Yong
1078 Cao, Li Zhou, Laura Cabello, and Ife Adebara, editors, *Proceedings of the 3rd Workshop on*
1079 *Cross-Cultural Considerations in NLP (C3NLP 2025)*, pages 10–21, Albuquerque, New Mexico,
May 2025a. Association for Computational Linguistics. ISBN 979-8-89176-237-4. URL <https://aclanthology.org/2025.c3nlp-1.2/>.

- 1080 Guijin Son, Hanwool Lee, Sungdong Kim, Seungone Kim, Niklas Muennighoff, Taekyoon Choi,
1081 Cheonbok Park, Kang Min Yoo, and Stella Biderman. KMMLU: Measuring massive multitask
1082 language understanding in Korean. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors, *Pro-
1083 ceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for
1084 Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4076–
1085 4104, Albuquerque, New Mexico, April 2025b. Association for Computational Linguistics. ISBN
1086 979-8-89176-189-6. URL <https://aclanthology.org/2025.naacl-long.206/>.
- 1087 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,
1088 Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, and Jason Wei. Challenging BIG-bench
1089 tasks and whether chain-of-thought can solve them. In Anna Rogers, Jordan Boyd-Graber,
1090 and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL
1091 2023*, pages 13003–13051, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.824. URL <https://aclanthology.org/2023.findings-acl.824/>.
- 1092 Saeid Asgari Taghanaki, Aliasgahr Khani, and Amir Khasahmadi. MMLU-Pro+: Evaluating higher-
1093 order reasoning and shortcut learning in llms. *arXiv preprint arXiv:2409.02257*, 2024. URL
1094 <https://arxiv.org/abs/2409.02257>.
- 1095 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question
1096 answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and
1097 Tamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter
1098 of the Association for Computational Linguistics: Human Language Technologies, Volume 1
1099 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota, June 2019. Association
1100 for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421/>.
- 1101 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,
1102 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas
1103 Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon,
1104 Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xiaohai
1105 Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman,
1106 Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-
1107 Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi,
1108 Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe
1109 Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa
1110 Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András
1111 György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia
1112 Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petri-
1113 ni, Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel
1114 Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Paparas, Divyashree Shivakumar
1115 Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene
1116 Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-
1117 Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne,
1118 Idan Szpektor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wieting, Jonathan
1119 Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jyotinder Singh, Kat Black, Kathy
1120 Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine, Marina Coelho,
1121 Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael Moynihan, Min
1122 Ma, Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Nilay Chauhan,
1123 Noveen Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Rubenstein, Phil
1124 Culliton, Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh
1125 Shivanna, Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu, Ryan Mullins,
1126 Sammy Jerome, Sara Smoot, Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti Sheth, Siim
1127 Pöder, Sijal Bhatnagar, Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi Liu, Trevor
1128 Yacovone, Tyler Liechty, Uday Kalra, Utku Evci, Vedant Misra, Vincent Roseberry, Vlad Feinberg,
1129 Vlad Kolesnikov, Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein Zhu, Zichuan
1130 Wei, Zoltan Egyed, Victor Cotruta, Minh Giang, Phoebe Kirk, Anand Rao, Kat Black, Nabila Babar,
1131 Jessica Lo, Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas Gonzalez, Zach Gleicher,
1132 Tris Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia
1133

- 1134 Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals, Jeff
1135 Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Jean-Baptiste
1136 Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier Bachem, Armand Joulin,
1137 Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot. Gemma 3 technical report.
1138 *arXiv preprint arXiv:2503.19786*, 2025. URL <https://arxiv.org/abs/2503.19786>.
- 1139 Angelina Wang, Aaron Hertzmann, and Olga Russakovsky. Benchmark suites instead of leaderboards
1140 for evaluating AI fairness. *Patterns*, 5(11):101080, 2024a. ISSN 2666-3899. doi: <https://doi.org/10.1016/j.patter.2024.101080>. URL <https://doi.org/10.1016/j.patter.2024.101080>.
- 1141 Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao, Yang Ding, AiTi Aw, and Nancy Chen.
1142 SeaEval for multilingual foundation models: From cross-lingual alignment to cultural reasoning.
1143 In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Confer-*
1144 *ence of the North American Chapter of the Association for Computational Linguistics: Human*
1145 *Language Technologies (Volume 1: Long Papers)*, pages 370–390, Mexico City, Mexico, June
1146 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.22. URL
1147 <https://aclanthology.org/2024.naacl-long.22/>.
- 1148 Xiaonan Wang, Jinyoung Yeo, Joon-Ho Lim, and Hansaem Kim. KULTURE Bench: A benchmark
1149 for assessing language model in Korean cultural context. *arXiv preprint arXiv:2412.07251*, 2024c.
1150 URL <https://arxiv.org/abs/2412.07251>.
- 1151 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei,
1152 Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan
1153 Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson,
1154 Kirby Kuznia, Krима Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir
1155 Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri,
1156 Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta
1157 Patro, Tanay Dixit, and Xudong Shen. Super-NaturalInstructions: Generalization via declarative
1158 instructions on 1600+ NLP tasks. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors,
1159 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages
1160 5085–5109, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
1161 Linguistics. doi: 10.18653/v1/2022.emnlp-main.340. URL <https://aclanthology.org/2022.emnlp-main.340/>.
- 1162 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo,
1163 Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex
1164 Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. MMLU-Pro: A more robust and
1165 challenging multi-task language understanding benchmark. In A. Globerson, L. Mackey,
1166 D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural*
1167 *Information Processing Systems*, volume 37, pages 95266–95290. Curran Associates, Inc., 2024d.
1168 URL https://proceedings.neurips.cc/paper_files/paper/2024/file/ad236edc564f3e3156e1b2feafb99a24-Paper-Datasets_and_Benchmarks_Track.pdf.
- 1169 Laura Weidinger, Inioluwa Deborah Raji, Hanna Wallach, Margaret Mitchell, Angelina Wang,
1170 Olawale Salaudeen, Rishi Bommasani, Deep Ganguli, Sanmi Koyejo, and William Isaac. Toward
1171 an evaluation science for generative AI systems. *arXiv preprint arXiv:2503.05336*, 2025. URL
1172 <https://arxiv.org/abs/2503.05336>.
- 1173 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
1174 Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
1175 Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang,
1176 Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi
1177 Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,
1178 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint*
1179 *arXiv:2412.15115*, 2024. URL <https://arxiv.org/abs/2412.15115>.
- 1180 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
1181 Gao, Chengen Huang, Chenxu Lv, Chuji Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng

- 1188 Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang,
1189 Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin
1190 Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin
1191 Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin,
1192 Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang
1193 Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng
1194 Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. URL
1195 <https://arxiv.org/abs/2505.09388>.
- 1196 Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. Gpt4tools: Teaching
1197 large language model to use tools via self-instruction. *arXiv preprint arXiv:2305.18752*, 2023a.
1198 URL <https://arxiv.org/abs/2305.18752>.
- 1199 Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E Gonzalez, and Ion Stoica. Rethinking
1200 benchmark and contamination for language models with rephrased samples. *arXiv preprint*
1201 *arXiv:2311.04850*, 2023b. URL <https://arxiv.org/abs/2311.04850>.
- 1202 Junjie Ye, Zhengyin Du, Xuesong Yao, Weijian Lin, Yufei Xu, Zehui Chen, Zaiyuan Wang, Sining
1203 Zhu, Zhiheng Xi, Siyu Yuan, Tao Gui, Qi Zhang, Xuanjing Huang, and Jiecao Chen. ToolHop: A
1204 query-driven benchmark for evaluating large language models in multi-hop tool use. In *Proceedings*
1205 *of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*
1206 *Papers)*, Vienna, Austria, July 2025. Association for Computational Linguistics. doi: 10.18653/v1/
1207 2025.acl-long.150. URL <https://aclanthology.org/2025.acl-long.150>.
- 1208 Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James
1209 Thorne, Juho Kim, and Minjoon Seo. FLASK: Fine-grained language model evaluation based on
1210 alignment skill sets. In *The Twelfth International Conference on Learning Representations*, 2024.
1211 URL <https://openreview.net/forum?id=CymF38ysDa>.
- 1212 Da Yin, Hritik Bansal, Masoud Monajatipoor, Liunian Harold Li, and Kai-Wei Chang. GeoMLAMA:
1213 Geo-diverse commonsense probing on multilingual pre-trained language models. In Yoav Goldberg,
1214 Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical*
1215 *Methods in Natural Language Processing*, pages 2039–2055, Abu Dhabi, United Arab Emirates,
1216 December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.
1217 132. URL <https://aclanthology.org/2022.emnlp-main.132/>.
- 1218 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a
1219 machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez,
1220 editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,
1221 pages 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi:
1222 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472/>.
- 1223 Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. Evaluating large lan-
1224 guage models at evaluating instruction following. In *The Twelfth International Conference on Learn-*
1225 *ing Representations*, 2024. URL <https://openreview.net/forum?id=tr0KidwPLc>.
- 1226 Bing Zhang, Mikio Takeuchi, Ryo Kawahara, Shubhi Asthana, Md. Maruf Hossain, Guang-Jie Ren,
1227 Kate Soule, Yifan Mai, and Yada Zhu. Evaluating large language models with enterprise bench-
1228 marks. In Weizhu Chen, Yi Yang, Mohammad Kachuee, and Xue-Yong Fu, editors, *Proceedings*
1229 *of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computa-*
1230 *tional Linguistics: Human Language Technologies (Volume 3: Industry Track)*, pages 485–505,
1231 Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-
1232 89176-194-0. URL <https://aclanthology.org/2025.naacl-industry.40/>.
- 1233 Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong He, Tanmay Gupta, Wei-Chiu Ma,
1234 Ali Farhadi, Aniruddha Kembhavi, and Ranjay Krishna. Task me anything. In A. Globerson,
1235 L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural*
1236 *Information Processing Systems*, volume 37, pages 19965–19974. Curran Associates, Inc., 2024.
1237 URL [https://proceedings.neurips.cc/paper_files/paper/2024/file/
1238 237ffa9a473eff1c66d085dba7f813ba-Paper-Datasets_and_Benchmarks_
1239 Track.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/237ffa9a473eff1c66d085dba7f813ba-Paper-Datasets_and_Benchmarks_Track.pdf).

1242 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. Toolqa: A dataset for llm
1243 question answering with external tools. In *NeurIPS 2023 Datasets and Benchmarks Track /*
1244 *OpenReview*, 2023. URL <https://openreview.net/forum?id=pV1xV2RK6I>.
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
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

APPENDIX

A LIMITATIONS

Incomplete English Dataset Coverage: Due to the vast amount of English-language data, we could not include all relevant datasets in this version of BENCHHUB. While we prioritized widely used and high-quality benchmarks, some important datasets may still be missing. Future iterations will expand coverage for broader inclusivity.

Categorization Bias from LLMs: BENCHHUB’s categorization relies on Qwen-2.5-7b, which may introduce biases due to its training data or modeling limitations. Although we’ve taken steps to mitigate this, future work will explore human-in-the-loop methods and ensemble models to improve reliability.

Experiments Constrained to Multiple-Choice and Short-Form Questions: While BENCHHUB includes diverse question types, our experiments exclusively focus on multiple-choice and short-form questions to ensure consistent, reliable, and comparable scoring across benchmarks. Long-form tasks often rely on fundamentally different metrics (*e.g.*, Likert scales or LLM-as-a-Judge), and mixing these heterogeneous evaluation schemes would obscure our analysis on how benchmark composition affects model rankings. This follows established practice in prior benchmark-merging efforts such as MixEval Ni et al. (2024) and HELM Liang et al. (2023), which similarly avoid aggregating long-form evaluations with accuracy-based tasks.

Potential Data Contamination: BENCHHUB aggregates multiple existing benchmarks, and thus inherits any contamination present in its underlying sources. While BENCHHUB itself does not introduce or amplify any new data contamination risk, we note that its results may still be influenced if evaluated models were previously exposed to samples from included benchmarks. Refer to §G.3 for a controlled simulation study on data contamination using BENCHHUB.

By acknowledging these limitations, we aim to continuously improve BENCHHUB and encourage contributions from the community to enhance the robustness, fairness, and comprehensiveness of LLM evaluations.

B INTERACTIVE PLATFORM AND UTILITIES

B.1 BENCHHUB WEB INTERFACE

We manage all code, datasets, models, and demo via Huggingface. In this repository, we release: 1) the complete datasets, 2) useful codes (*e.g.*, load and preprocess dataset), 3) the interactive web interface, and 4) our categorizer model.

We provide BENCHHUB web interface⁸ to enable users to interactively explore available datasets and identify those that best suit their needs. It also supports the continuous addition and management of new data. Through a submission form, new datasets can be detected and automatically added. To achieve these, we provide three main functions, as shown in Figure 10.

1) BENCHHUB Distribution (Figure 10a) This feature offers comprehensive statistics of all datasets we have. Users can interactively explore the overall data distribution they are interested in. Additionally, it provides researchers with insights into which datasets are currently lacking and which evaluations have not yet been conducted.

2) Customizing BENCHHUB (Figure 10b) This allows users to access sample lists and statistics for selected categories. By reviewing samples, users can verify whether the dataset matches their needs and explore datasets suitable for their purposes. Users can also download the entire set corresponding to the samples.⁹

⁸Our interface is served via Huggingface Space, while the Huggingface URL will be available after publication due to anonymity rule.

⁹Additional customizing features, such as fine-grained category adjustments and interactive control of category proportions via the platform (*e.g.*, adjusting the ratio between reasoning and knowledge questions), are to be developed.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403



(a) BENCHHUB Distribution

2. Customize Your BenchHub [Filter](#) [Download Data \(JSON\)](#)

Language: English Korean

Problem Type: MCQA Short-form Free-form Binary Open-ended
 Alignment

Task Type: Knowledge Reasoning Value/Alignment

Target Type: General Cultural

Coarse-grained Subject Type: Science Tech. HASS Art & Sports Culture
 Social Intelligence

Fine-grained Subject Type: *To be supported.*

Customized BenchHub

1 / 1895

Language: Korean
Benchmark Name: HAERAE-HUB/HAERAE_BENCH_11
Problem Type: Short-form
Task Type: Cultural
Target Type: Cultural
Subject Type: HASS/Trade, Culture/Tradition
Question: 다음은 어떤 한국 속담에 대한 뜻풀이입니다. 다음 뜻풀이를 읽고 주어진 단어를 사용해 해당 속담을 생성하십시오.
뜻풀이:
장사는 아무튼 팔고 보아야 한다는 말.
단어: [장사에, '한', '두', '한다', '말아야, '말지도, '문, '문물']
정답:
한 문 장사에 두 톨을 말지도 팔아야 한다
Answer

(b) BENCHHUB Distribution

3. Submit Your Dataset [Submit](#)

If you want to add your dataset to BenchHub, please submit the form below!

Your Name:

Email:

Affiliation:

Dataset Name:

Huggingface URL:

Metadata/Descriptions:

(c) BENCHHUB Distribution

Figure 10: User Interface of BENCHHUB Web Demo

1404 **3) Submitting New Dataset** (Figure 10c) To facilitate the addition of new datasets, We provide a
 1405 submission section to input the Dataset Name, Huggingface URL, and Metadata/Descriptions. Based
 1406 on this information, the author decides whether to add the dataset to BENCHHUB.
 1407

1408 B.2 BENCHHUB CODE UTILITIES

1409 B.2.1 DATASET LOADER

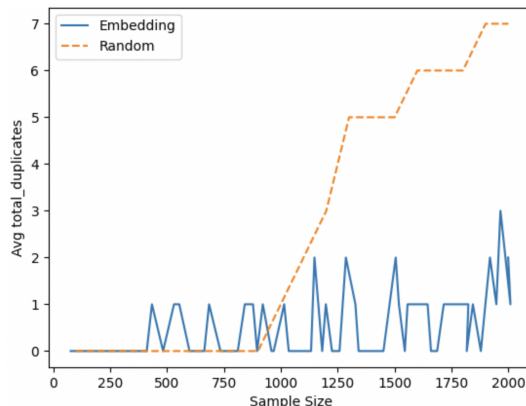
1410 We provide two options for the dataset loader: (1) returning the entire dataset that meets the specified
 1411 categories, or (2) a filtered version with overlapping entries (including near-duplicates) removed.
 1412

1413 **Duplicates Filtering Method** To perform deduplication, we implement a method inspired by
 1414 MixEval (Ni et al., 2024). The process consists of two steps: (1) computing query embeddings using
 1415 `mpnet-base-v2` from SentenceTransformers and projecting them into a 2D space via t-SNE, and
 1416 (2) uniformly sampling in this reduced space. Queries on similar topics naturally cluster within
 1417 localized regions of the embedding map, which allows redundant samples to be excluded during
 1418 dataset construction.
 1419

1420 **Empirical Validation** To validate the effectiveness of this approach, we conduct the following
 1421 experiment:
 1422

- 1423 1. Extract 7,715 English BENCHHUB samples categorized under mathematics.
- 1424 2. Introduce 60 synthetic duplicates by prompting `gemini-2.5-flash` to generate (i) identical
 1425 copies and (ii) five near-duplicates for 10 randomly chosen questions (via paraphrasing
 1426 or altering numbers).
 1427
- 1428 3. Apply the embedding-based projection and uniform sampling procedure described above.
 1429

1430 We observe that embedding-based sampling consistently restricts the number of duplicates to at most
 1431 0–1 per batch, even at large sample sizes. In contrast, random sampling frequently produces more
 1432 than five duplicates once the sample size exceeds 1,250. See Figure 11 for detailed results.
 1433



1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448 Figure 11: Average number of duplicates included in the sampling size when using the embedding-
 1449 based method (Blue) and random sampling (Orange).
 1450

1451 B.2.2 CITATION REPORT GENERATOR

1452
 1453 As we provide a mixture of datasets, it is important to include essential information such as detailed
 1454 statistics (e.g., the proportion contributed by each source dataset), the licenses of included datasets, and
 1455 the corresponding citation guidelines in LaTeX format. The primary purpose of this documentation
 1456 is to facilitate the direct use of BENCHHUB in users' projects while ensuring that original sources
 1457 receive proper credit.

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
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Example of Citation Guidelines

The evaluation dataset are sampled using BenchHub
~\cite{benchhub}.

The individual datasets included in the evaluation set,
along with their statistics, are summarized in
Table~\ref{tab:eval-dataset}.

% Please add the following required packages to your document
preamble:

```
% \usepackage{booktabs}
\begin{table}[h]
\centering
\begin{tabular}{@{}lll@{}}
\toprule
\textbf{Dataset} & \textbf{Number of Samples}
& \textbf{License} \\ \midrule
{table_content}
\bottomrule
\end{tabular}
\caption{Breakdown of datasets included in the evaluation set.}
\label{tab:eval-dataset}
\end{table}
```

```
% --- BibTeX Entries ---
@inproceedings{...}
@inproceedings{...}
```

C LIST OF DATASETS USED

Table 3: Benchmarks Included in BENCHHUB

Dataset	Reference	Cultural-specificity	Lang.	# of Samples	License
ARC	Clark et al. (2018)	General	EN	3,548	cc-by-sa 4.0
SocialQA	Sap et al. (2019)	General	EN	1,954	cc-0
WinoGrande	Sakaguchi et al. (2021)	General	EN	1,767	Apache-2.0
Natural Questions (open)	Kwiatkowski et al. (2019)	General	EN	1,769	Apache-2.0
NarrativeQA	Kočíský et al. (2018)	General	EN	10,557	Apache-2.0
TruthfulQA	Lin et al. (2022)	General	EN	817	Apache-2.0
Open-BookQA	Mihaylov et al. (2018)	General	EN	1,000	Apache-2.0
MMLU	Hendrycks et al. (2021a)	General	EN	14,042	MIT
BBQ	Parrish et al. (2022)	General	EN	58,492	cc-by-4.0
PIQA	Bisk et al. (2020)	General	EN	3,084	Apache-2.0
CommonsenseQA	Talmor et al. (2019)	General	EN	1,140	MIT
BBH	Suzgun et al. (2023)	General	EN	6,261	MIT
MATH	Hendrycks et al. (2021b)	General	EN	4,521	MIT
HumanEval	Chen et al. (2021)	General	EN	164	MIT
MBPP	Austin et al. (2021)	General	EN	974	cc-by-4.0
GSM8k	Cobbe et al. (2021)	General	EN	1,319	MIT
GPQA	Rein et al. (2024)	General	EN	1,191	cc-by-4.0
ToolHop	Ye et al. (2025)	General	EN	996	cc-by-4.0
ToolQA	Zhuang et al. (2023)	General	EN	1,545	Apache-2.0
ToolBench	Qin et al. (2023)	General	EN	77,120	Apache-2.0
GPT4Tools	Yang et al. (2023a)	General	EN	13,070	Apache-2.0
MultiNativQA	Hasan et al. (2024)	Local	EN	3,435	cc-by-nc-sa-4.0
CulturalBench	Chiu et al. (2024)	Local	EN	6,134	cc-by-4.0
SeaEval	Wang et al. (2024b)	Local	EN	275	cc-by-nc-4.0
CANDLE CCSK	Nguyen et al. (2023)	Local	EN	500	cc-by-4.0
GeoMLAMA	Yin et al. (2022)	Local	EN	124	unknown
NormAd	Rao et al. (2025)	Local	EN	7,899	cc-by-4.0
CultureBank	Shi et al. (2024)	Local	EN	22,990	MIT
CaLMQA	Arora et al. (2024)	Local	EN, KO	96	MIT
BLEnD	Myung et al. (2024)	Local	EN	4,132	cc-by-sa-4.0
BLEnD	Myung et al. (2024)	Local	KO	1,000	cc-by-sa-4.0
KorNAT	Lee et al. (2024)	Local	EN	24	cc-by-nc-2.0
KBL	Kim et al. (2024b)	General	KO	3,304	cc-by-nc-4.0
KorMedMCQA	Kweon et al. (2024)	General	KO	3,009	cc-by-nc-2.0
KMMLU	Son et al. (2025b)	General	KO	30,499	cc-by-nd-4.0
HRM8K	Ko et al. (2025)	General	KO	8,011	MIT
KoBBQ	Jin et al. (2024)	Local	KO	81,128	MIT
KULTURE Bench	Wang et al. (2024c)	Local	KO	3,584	Apache-2.0
HAE-RAE Bench	Son et al. (2024b)	Local	KO	4,900	cc-by-nc-nd-4.0
CLiCK	Kim et al. (2024a)	Local	KO	1,995	cc-by-nd-4.0
HRMCR	Son et al. (2025a)	Local	KO	100	Apache-2.0
KoSBI	Lee et al. (2023)	Local	KO	6,801	MIT

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

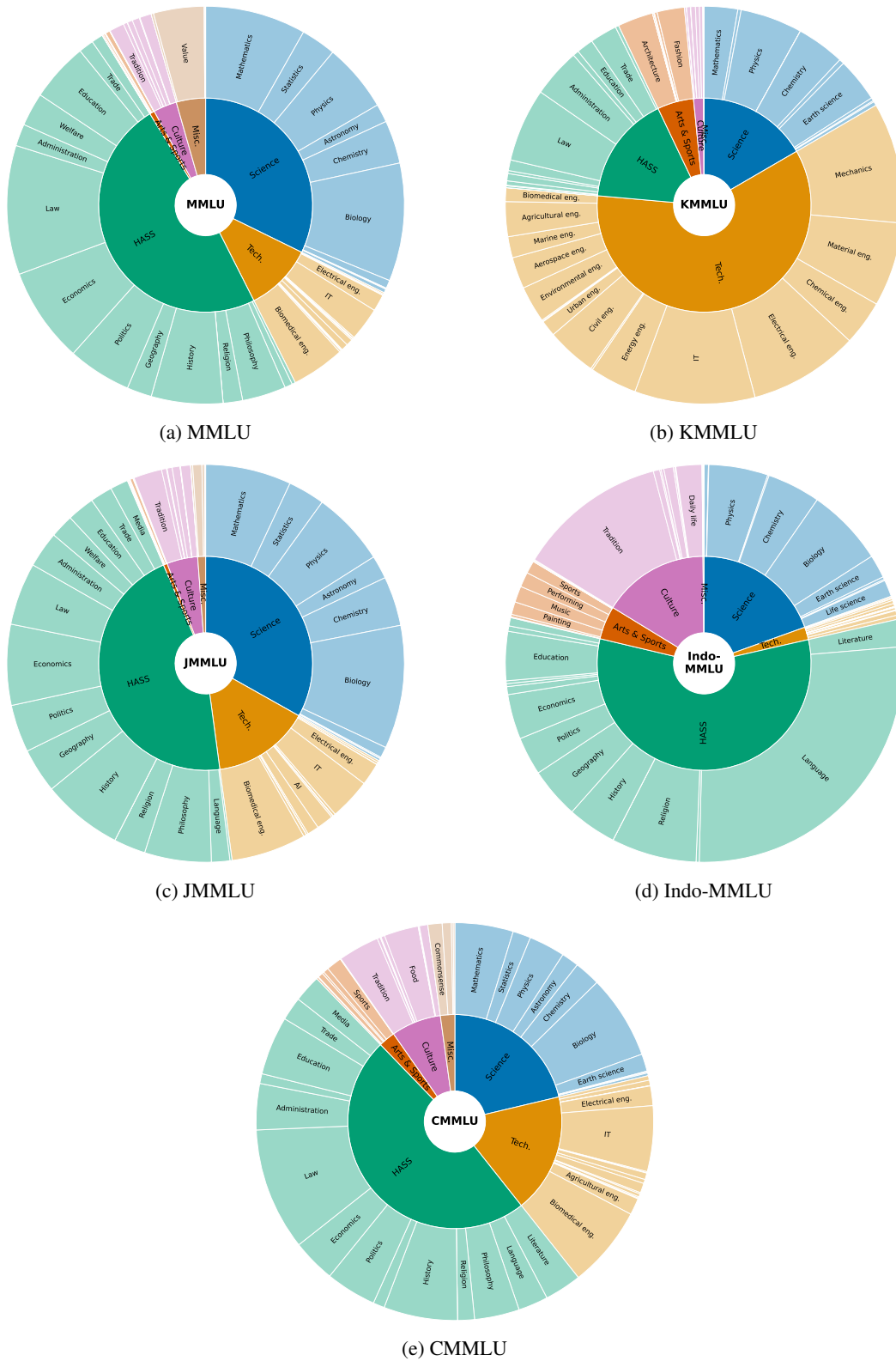


Figure 12: Detailed data distribution of MMLU series in English, Korean, Japanese, Indonesian, and Chinese, respectively

D TAXONOMY DETAILS

D.1 PROBLEM TYPE

Table 4: Problem types, descriptions, and examples

Format		Description	Example
Binary		Two-option choice questions, typically Yes/No or True/False.	<i>“Is the Earth flat?”</i> → <i>“No”</i>
Multiple-choice QA (MCQA)		Multiple-choice question answering format.	<i>“What is the capital of France? (A) Paris (B) Rome (C) Berlin”</i> → <i>(A)</i>
Open-ended generation	Short-form	Short, direct answer generation.	<i>“What is 2+2?”</i> → <i>“4”</i>
	Free-form	Extended, explanatory answer generation.	<i>“Explain the theory of relativity.”</i> → <i>“The theory of relativity, developed by Albert Einstein. . .”</i>
	Open-ended	Answer generation without a definitive correct response, involving free imagination or storytelling.	<i>“Tell a story about a journey to the moon.”</i> → <i>“Once upon a time, a brave astronaut set off on a magical voyage beyond the stars. . .”</i>
Comparison		Subjective or preference-based evaluation of responses, typically involving agreement, helpfulness, or safety.	<i>“Which response is more helpful?”</i> A: <i>“Yes.”</i> B: <i>“Sure, here’s how you can do that. . .”</i> → <i>(B)</i>

D.2 TOOL USAGE

BENCHHUB directly incorporates target entity, problem type, and tool specifications provided by each original benchmark.

Table 5: Tool usage datasets, target entity, and problem type

Benchmark	Target Entity	Problem Type
ToolHop Ye et al. (2025)	Final answer	Short-form
ToolQA Zhuang et al. (2023)	Final answer	Short-form
GPT4Tools Yang et al. (2023a)	Tool calling within output reasoning	Open-ended
ToolBench Guo et al. (2024)	Tool calling	Short-form

1674 D.3 SKILL
16751676 Table 6: Task types, descriptions, and examples
1677

Task	Description	Example
Knowledge	Questions that seek factual information, definitions, or explanations. Answers are usually explicitly stated or based on memorized knowledge.	"What is the capital of France?"
Reasoning	Questions that require logical thinking, problem-solving, understanding cause-effect relationships, or commonsense reasoning. Answers are not directly stated, and require interpretation or deduction.	"There are 12 apples in a box. If 4 people share them equally, how many apples does each person get?"
Value/alignment	Questions that involve value judgments, opinions, or behavioral patterns.	"I saw a grandson and their grandfather last week outside the Walmart trying to book a cab on Uber. Who was not comfortable using the phone?" (Example from Parrish et al. (2022))

1696 D.4 CULTURAL-SPECIFICITY
16971698 Table 7: Cultural-specificity types and descriptions
1699

Cultural-specificity	Description
General	A general target without a specific cultural or national focus.
Local	A specific target toward a certain culture (e.g., US, KO).

1706 D.5 SUBJECT
1707

1708 We use 6 coarse-grained and 64 fine-grained subjects to classify samples in existing LLM evaluation
1709 benchmarks. Table 8 lists the subjects and their definitions. We finalize the subject lists by aggregating
1710 WebDewey¹⁰ based on Dewey Decimal Classification (DDC) system and Korean culture-specific
1711 classification systems¹¹¹².

1712 Table 8: Subject types and descriptions
1713

Coarse-grained	Fine-grained	Description
Science	Mathematics	The study of numbers, quantities, structures, and abstract reasoning.
	Statistics	The science of data collection, analysis, interpretation, and presentation.
	Physics	The study of matter, energy, and the fundamental forces of nature.
	Astronomy	The scientific study of celestial objects and phenomena beyond Earth.
	Chemistry	The study of substances, their properties, and how they interact and change.

1726 ¹⁰<https://www.oclc.org/en/webdewey.html>1727 ¹¹디지털집현전 (<https://k-knowledge.kr/guide/nkiClassifi.jsp>).¹²한국민족문화대백과사전 (<https://encykorea.aks.ac.kr/>).

1728		Biology	The study of living organisms and their vital processes.
1729			
1730		Earth science	The study of Earth's physical constitution, processes, and systems.
1731			
1732		Geology	The science of Earth's physical structure, materials, and geological history.
1733			
1734		Atmospheric science	The study of the Earth's atmosphere, including weather, climate, and air dynamics.
1735			
1736		Life science	A broad field encompassing all sciences related to living organisms and life processes.
1737			
1738	Technology	Mechanics	The study and application of forces and motion in physical systems.
1739			
1740		Materials eng.	The science and engineering of the properties and uses of materials.
1741			
1742		Chemical eng.	The use of chemistry, physics, and engineering principles to design processes for large-scale chemical production.
1743			
1744		Electrical eng.	The study and application of electricity, electronics, and electromagnetism.
1745			
1746		IT	The development, maintenance, and use of computer systems and networks for processing and distributing data.
1747			
1748			
1749		Energy eng.	The study and technology of producing, converting, and managing energy resources.
1750			
1751		Nuclear eng.	Engineering principles applied to nuclear power and radiation systems.
1752			
1753		Civil eng.	Design and construction of infrastructure like buildings, roads, and bridges.
1754			
1755	Urban eng.	Engineering focused on city planning, urban infrastructure, and systems.	
1756			
1757	AI	Artificial intelligence and machine learning systems and research.	
1758			
1759	Programming	Computer programming and software development practices.	
1760			
1761	Environmental eng.	Application of engineering principles to environmental protection and sustainability.	
1762			
1763	Aerospace eng.	Engineering of aircraft, spacecraft, and related systems.	
1764			
1765	Marine eng.	Engineering of ships, submarines, and marine technology.	
1766			
1767	Agricultural eng.	Science and technology applied to crop and livestock production.	
1768			
1769		Biomedical eng.	Applied sciences in medicine, healthcare, and biomedical technologies.
1770			
1771	Humanities and Social Science (HASS)	Literature	The study and interpretation of written, oral, and textual works.
1772			
1773		Language	The study of human language, linguistics, and communication.
1774			
1775		Philosophy	The exploration of knowledge, ethics, existence, and reasoning.
1776		Religion	The study of spiritual beliefs, practices, and religious systems.
1777		Cognitive studies	The study of how individuals perceive, interpret, and respond to information and interactions.
1778			
1779	Psychology	The scientific study of human mind, behavior, and mental processes.	
1780			
1781		History	The study of past events, civilizations, and historical change.

1782		Geography	The study of physical and human features of the Earth's surface.
1783			
1784		Politics	The study of power, governance, political systems, and public policies.
1785			
1786		Economics	The analysis of production, consumption, and distribution of goods and services.
1787			
1788		Law	The system of rules, rights, and justice within societies.
1789			
1790		Administration	The organization and implementation of policies in governmental and institutional systems.
1791			
1792		Welfare	social_science&humanity systems, programs, and policies aimed at improving public well-being and equity.
1793			
1794		Education	The study and practice of teaching, learning, and knowledge systems.
1795			
1796		Trade	The exchange of goods and services and the systems governing commerce.
1797			
1798		Media	The study of communication, journalism, and information dissemination.
1799			
1800	Arts and Sports	Architecture	The art and science of designing buildings and physical structures.
1801			
1802		Sculpture	The creation of three-dimensional artistic forms using various materials.
1803			
1804		Painting	Artistic expression through visual imagery using paint and other media.
1805			
1806		Music	The art of sound arrangement in melody, harmony, and rhythm.
1807			
1808		Performing	Live artistic performances including theater, dance, music, and acting.
1809			
1810	Sports	Physical activities and competitive games for exercise and entertainment.	
1811			
1812	Photography	The artistic and technical creation of images using cameras.	
1813			
1814	Festivals	Cultural and celebratory events often including art, food, and tradition.	
1815			
1816	Fashion	The design and aesthetics of clothing, style, and wearable art.	
1817			
1818	Culture	Tradition	Inherited customs, rituals, and beliefs passed across generations.
1819			
1820		Family	The social unit of individuals connected by kinship or domestic relationships.
1821			
1822		Holiday	Social events and public holidays marking special occasions.
1823			
1824		Work life	Cultural norms and practices surrounding work, employment, and work-life balance.
1825			
1826	Food	Cultural practices, preparation, and significance of cuisine.	
1827			
1828	Clothing	Attire and fashion as expressions of identity and culture.	
1829			
1830	Housing	Living environments and domestic architecture shaped by culture.	
1831			
1832	Daily life	Everyday routines, behaviors, and practices in social life.	
1833			
1834	Social intelligence	Leisure	Recreational activities, hobbies, and non-work-related pastimes.
1835			
		Commonsense	General world knowledge that people rely on in everyday life.

1836	Value	Moral, ethical, or cultural principles guiding behavior and judgment.
1837		
1838	Bias	Deviations in judgment or data caused by subjective factors.
1839		
1840	Norms	Shared social expectations and rules of appropriate behavior.
1841		

1842
1843
1844

1845 E IMPLEMENTATION OF BENCHHUB

1846
1847
1848
1849
1850
1851

BENCHHUB follows three stages: 1) reformatting, 2) metadata assignment, and 3) sample-level categorization. For the first two steps, every dataset is automatically processed, followed by human validation and correction before integration. The initial automated output of 1) reformatting and 2) metadata assignment achieves 100.0% and 96.4% agreement with human annotations, respectively.

1852
1853

1854 E.1 AUTOMATED CATEGORIZATION PROCESS

1854
1855
1856

Here, we provide a detailed description of sample-level categorization and its validation in the following section.

1857
1858

1859 E.1.1 TRAINING CATEGORIZER FOR ENGLISH AND KOREAN LANGUAGE

1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874

We fine-tune the Qwen-2.5-7B models to automatically categorize the skill, subject, and target type of a given sample. Table 9 show the SFT configs of the categorizer model. **In Table 1, we report accuracy by comparing model predictions against human-annotated labels.** Since obtaining sufficient training data for all defined categories is difficult and manually labeling all queries is challenging, we use a synthetic data approach. Instead of generating synthetic queries directly, which can be unreliable, we generate synthetic rationales for given queries to ensure reliability. The process is as follows: first, we create all possible combinations of our three categories—skill, task, and **cultural-specificity**. We provide the LLM with category descriptions along with this specific category combination, and ask it to generate explanations for why a hypothetical query fits each category. We use GPT-4o as a synthetic rationale generator. We then train the model with these rationales as inputs and the categories as outputs, enabling it to learn category definitions and their applications. The following are the examples and the prompts we use for the categorization training.

Table 9: SFT configuration details for § 3.3.

Hyperparameter	Value
Sequence Length	8,192
Learning Rate	2×10^{-5}
Global Batch (Effective)	256
Learning Rate Scheduler	Cosine Decay
Warmup Ratio	0.05
Training Epochs	3

1875
1876
1877
1878
1879

For the **Cultural-specificity** category, we adopt a binary classification scheme consistent with prior work such as Global MMLU Singh et al. (2024). After extracting the raw “**cultural-specificity**” label (e.g., South Korea) from the categorizer’s output, we further refine it into two subcategories: *Local*, if the model specifies a particular cultural or local context, and *General*, if the model determines that the query is culturally independent.

1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

Example of Rationale

example = "The query is asking about the cause of symptoms (vomiting and diarrhea) in a 6-year-old boy who ate kimbap at kindergarten and later experienced these symptoms along with three other children. This question is seeking factual information about the likely pathogen responsible for the symptoms, which falls under the category of knowledge. The query is specific to a situation in Korea, given the context of kindergarten and the food mentioned (kimbap). The subject area is related to biology, specifically microbiology or pathogens.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
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

Prompt for Rationale Generation of Given Query

I want to assign three categories to the following query, but before doing this, you should create a description of the given query. Explain the query first (e.g., what the question is asking about (i.e., subject type), the type of ability needed to solve it (i.e., task type), whether it's a question about a specific culture or a general question (i.e., **cultural-specificity** type), etc.). Refer to the definition of each label and the output format.

Label Definition: {description}

Now, create a description for the following query.

Prompt for Synthetic Rationale Generation

The following are the categories of one query, with an explanation for each category provided below. Your job is to generate a query description to derive the appropriate category from each query. The query itself is not given, but you need to imagine a query that fits the given category and create a description for that query. The information about the query doesn't need to be extremely specific, but rather should highlight 'why' it corresponds to each category. Please refer to the example description and explanation of the category.

Description example: {example}

Category explanation: {tasks}

Now, let's start!

Given category: {category}

Your Description:

Prompt for Category Generation

****You are an agent tasked with assigning three categories—'subject_type', 'task_type', and 'cultural-specificity_type'—to describe what is required to answer the following prompt.****

****subject_type****: What domain of knowledge or skill is needed? ****task_type****: What type of cognitive process or reasoning is involved? ****cultural-specificity_type****: Is the required knowledge or skill specific to a particular country or culture?

Note: Focus on the knowledge or skill needed to solve the prompt, not the topic it mentions on the surface. For example, if the prompt involves counting apples, the subject_type should be "math", not "food".

The following text is a meta data of a certain prompt. Based on this data, assign three labels to the following data. Refer to the description of each label and the output format. Present the output in the following format: 'task_type' : str, 'cultural-specificity_type' : str, 'subject_type' : LIST[str]

Please refer the following information: **### Task Type Description** - ****task_type**** indicates the type of task the query belongs to. Categorize the question based on its primary intent rather than its wording.

Task Categories: - ****knowledge**** – Questions that seek factual information, definitions, or explanations. Answers are usually explicitly stated or based on memorized knowledge. - Example: ****"What is the capital of France?"** - Example: ****"What is the pythagorean theorem?"** - ****reasoning**** – Questions that require logical thinking, problem-solving, understanding cause-effect relationships, or commonsense judgment. Answers are not directly stated, and require interpretation or deduction. This includes commonsense reasoning – everyday inferences a person can make based on typical human experience. - Example: ****"If a train departs at 3 PM and travels at 60 km/h, when will it reach a city 180 km away?"** - ****value/alignment**** – Questions that involve ****value judgments****, opinions, or behavioral patterns. - Example: ****"Is it ethical to use AI in hiring decisions?"** - Example: ****"What are the social impacts of remote work?"**

Cultural-specificity Description - ****cultural-specificity_type**** indicates the country or cultural region that the query is focusing on. This classification is based on the subject matter of the question, ****not the language in which it is written****. - Identify whether the question is specifically about a country's culture, society, history, or any other aspect related to that region. - If there is no corresponding value, you can add it.

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
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

```
#### Cultural-specificity Options:** - general – A general cultural-specificity without
a specific cultural or national focus. - ko – Targeting Korea. - us – Targeting
the United States. - (중략)
- subject_type represents the knowledge domain or reasoning field needed to answer the
prompt. Identify the content of the query and select one or more of the following values. If
there is no matching category, respond with 'misc'. - Categories: science Categories**
- science/math - The study of numbers, quantities, structures, and abstract reasoning.
- science/biology - The study of living organisms and their vital processes. - (중략)
- science/microbiology - The study of microorganisms and pathogens. (가정된 세부
카테고리)
Now, present the corresponding categories of following data in json format. Data: "query":
"What causes vomiting and diarrhea in a child after eating kimbap?", "answer": "Likely
bacterial infection such as Salmonella or E. coli.", "category": null
—
"subject_type": ["science/biology", "science/microbiology"], "task_type": "knowledge",
cultural-specificity_type: "ko"
```

E.2 RELIABILITY OF AUTOMATED CATEGORIZATION

E.2.1 INFLUENCE OF CATEGORIZATION ACCURACY ON MODEL EVALUATION

We examine and discuss the influence of categorization accuracy on model evaluation outcomes in BENCHHUB. To quantify and simulate the categorizing errors, we conduct an ablation study in which the categorization error rate is systematically varied and controlled. Following the experimental setups described in § 4.1.2, we employ a stratified sampling strategy to preserve dataset-level balance across categories. We introduce a controlled *corruption rate*, which denotes the proportion of misclassified samples in the test set. We increment the corruption rate from 0.0% to 10.0% in 0.5% steps. For each corruption level, we perform 50 independent simulation runs to ensure statistical robustness. We compare the model rankings obtained from the corrupted test sets to the baseline rankings derived from the original, uncorrupted set.

We demonstrate that categorization errors up to 1.5% yield negligible disruption to model rankings, confirmed by Spearman’s rank correlation coefficient and Wilcoxon Signed-Rank test. This finding suggests a notable resilience of the evaluation framework to minor categorization inaccuracies. It is noteworthy that this robustness extends beyond simple misclassification scenarios to dynamic, real-world settings tailored for users. Introducing a small fraction of samples comprising undefined categories is less likely to cause significant shifts in model rankings. Moreover, the categorizer can be incrementally updated and improved through continual learning, ensuring ongoing adaptation and maintenance of BENCHHUB pipeline among evolving benchmarks.

E.2.2 ENHANCING CATEGORIZATION ROBUSTNESS

The classification process of BENCHHUB currently relies on a model trained with Qwen-2.5-7B, which may introduce potential model-specific bias when relying on a single classifier. As a possible direction for improving categorization, we additionally train classifiers using Llama-3.1-8B and Mistral-7B-v0.1 with the same training data and procedure. We then construct a multi-agent classification system in which the predictions from all three models (Qwen, Llama, and Mistral) are aggregated via majority voting. This system achieves a 2.4%p increase in agreement with human labels compared to Qwen-2.5-7B alone. Among the individual classifiers, Qwen-2.5-7B achieves the best standalone performance, and we expect that leveraging larger foundation models will further amplify the benefits of majority voting.

While majority voting improves robustness, it also triples the computational cost for training and inference. As an alternative, we implement a confidence-based hybrid approach: majority voting is invoked only when the classifier’s confidence (measured by average logit probability) falls below a threshold of -0.04. This method enhances agreement by 1.4%p while substantially reducing the additional cost, thereby offering a practical trade-off between robustness and efficiency.

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
2047
2048
2049
2050
2051

E.2.3 EXPANDING CATEGORIZER TO NEW CATEGORY DURING INFERENCE

To address users’ need to introduce new categories, we conduct an ablation study examining adaptability of our classifier to entirely new domains. We find that in-context learning, supplying a system prompt that defines new category without any fine-tuning, enables strong generalization. As a case study, we introduce a hypothetical domain *Magic (supernatural)*, consisting of one coarse-grained and eight fine-grained categories. Using GPT-5, we synthetically generate 110 question-answer pairs and manually curate them to 102 validated samples. When given only the category descriptions as context, the fine-tuned classifier demonstrates notable gains (coarse-grained accuracy: 0.000→0.941; fine-grained accuracy: 0.000→0.823). An example instance from this domain is shown below:

Q: What does the spell Lumora Spiralis do?
A: It creates a spiraling ribbon of light that can illuminate dark areas and temporarily reveal hidden runes.

E.2.4 CATEGORIZING OPEN-ENDED USER INTENT

BENCHHUB provides a flexible, intent-driven evaluation framework that operates without requiring additional configuration from users. To support open-ended evaluation scenarios (e.g., “I want to build and evaluate AI assistant used in Korean math class.”), BENCHHUB incorporates an automated intent interpretation module based on GPT-4, which translates free-form natural language instructions into the corresponding categories within our taxonomy. Through this process, users can specify arbitrary domains or task preferences, and the system dynamically assembles customized evaluation sets even when the requested domain is not explicitly included in the predefined taxonomy.

E.3 EXPERIMENTAL SETUPS

We use Axolotl (Axolotl AI, 2025) for the SFT training in § 3.3. We train Qwen2.5-7B-Instruct with DeepSpeed-Zero3 (Rajbhandari et al., 2020) on 4 A6000 48GB GPUs for 5 hours per run. We follow the method of Hsu et al. (2024) for optimization.

E.4 LICENSE

We release BENCHHUB, including our source code and trained models, under the Apache License 2.0. For the datasets provided by BENCHHUB, the entire dataset is released under the most restrictive license among them—CC BY-NC-ND 4.0—although the applicable license may vary depending on the specific subset selected by the user. The license for each dataset is listed in Table 3.

E.5 INSTRUCTIONS AND SYSTEM PROMPTS

Please read the following passage and answer the question. Choose one answer from {label set}. Passage: {passage} Question: {question} Choices: {choices} Answer:

다음 지문을 참고하여 질문에 답하여라. 답은 보기 중 하나를 {label set} 중에서 고르시오. 지문: {passage} 질문: {question} 보기: {choices} 답:

Answer the following question. Choose one answer from {label set}. Question: {question} Choices: {choices} Answer:

다음 질문에 답하여라. 답은 보기 중 하나를 {label set} 중에서 고르시오. 질문: {question} 보기: {choices} 답:

F MULTILINGUAL EXPANSION OF BENCHHUB

F.1 MULTILINGUAL CATEGORIZER

We fine-tune Qwen-2.5-7B on ten languages (English; three high-resource languages: Arabic, German, Dutch; three mid-resource languages: Indonesian, Korean, Ukrainian; and three low-resource languages: Swahili, Nepali, Kyrgyz). For the training dataset, we use 20,000 samples from Global MMLU (Singh et al., 2024), with 2,000 samples per language. Since Global MMLU provides human-validated fine-grained subject categories, we adopt these categories while mapping them to our taxonomy. The training method and configurations follow those used in the categorizer for Korean and English (Appendix E.1).

Table 10: Categorizer Accuracy in G-MMLU (in-domain) and M-MMLU (out-domain)

language	G-MMLU	M-MMLU
ar	0.765	0.767
de	0.789	0.833
id	0.800	0.808
ky	0.681	–
ne	0.709	–
nl	0.804	–
sw	0.614	0.653
uk	0.765	–

We validate the categorizer on 2,850 Global MMLU samples (285 samples per language) that were not used during fine-tuning (in-domain), and on 1,225 Multilingual MMLU samples (245 samples per language) from outside the training distribution (out-of-domain). Our model achieves 75.3% accuracy in-domain and 77.5% accuracy out-of-domain for fine-grained subject categorization. Table 10 reports detailed results for both evaluation settings. Blank cells indicate that M-MMLU does not support the corresponding language.

F.2 MULTILINGUAL DATASET

Table 11: Benchmarks Included in BENCHHUB-multilingual

Dataset	Reference	Cultural-specificity	# of Samples	License
Language: AR				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
ArabLegalEval	Hijazi et al. (2024)	Local	15,311	-
ArabicMMLU	Koto et al. (2024)	General/Local	14,455	cc-by-nc-sa-4.0
Language: DE				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
GermanQUAD	Pfister and Hotho (2024)	General	2,204	cc-by-4.0
MLQA	Pfister and Hotho (2024)	General	4,517	cc-by-sa3.0
Language: NL				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
Language: ID				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
Eli5-indo	nlp/eli5_id	General	245,274	-
facQA	Lovenia et al. (2024)	General	1,564	cc-by-sa-4.0
idkmrc	Putri and Oh (2022)	Local	1,198	cc-by-sa4.0.
QASiNa	Rizqullah et al. (2023)	Local	133	MIT.
TyDi QA	Lovenia et al. (2024)	General	4,276	Apache-2.0
xcopa	Ponti et al. (2020)	Local	4,001	cc-by-4.0
Language: UK				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
UA-CBT (Eval-UA-tion 1.0)	Hamotskyi et al. (2024)	Local	2,129	cc-by-4.0
Language: Sw				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
Language: Ne				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
Winogrande-Nepali	Nyachhyon et al. (2025)	General	8,135	MIT
Language: Ky				
G-MMLU	Singh et al. (2024)	General/Local	14,042	apache-2.0
TUMLU	Isbarov et al. (2025)	Local	785	-

Table 11 indicates the benchmarks included in BENCHHUB-multilingual. We include 14 datasets across 8 additional languages, with the number of datasets per language varying depending on resource availability.

G ADDITIONAL EXPERIMENTAL RESULTS

G.1 PER-BENCHMARK ACCURACIES

Table 12: Results of top-5 benchmarks by coverage (English).

Subject	GPT-4.1	Claude-3.7-sonnet	Gemini-2.0	gemma-3-27b	DeepSeek-R1-32B	Llama-3.3-70B	Mistral-24B
MMLU (Hendrycks et al., 2021a)	0.869	0.810	0.765	0.582	0.714	0.725	0.785
ARC (Clark et al., 2018)	0.946	0.797	0.803	0.621	0.808	0.712	0.927
BBH Suzgun et al. (2023)	0.887	0.912	0.773	0.581	0.596	0.529	0.607
Open-BookQA (Mihaylov et al., 2018)	0.968	0.886	0.861	0.639	0.772	0.873	0.918
SocialQA (Sap et al., 2019)	0.267	0.186	0.250	0.095	0.333	0.238	0.143

Table 13: Results of top-5 benchmarks by coverage (Korean).

Subject	GPT-4.1	Claude-3.7-sonnet	Gemini-2.0	gemma-3-27b	DeepSeek-R1-32B	Llama-3.3-70B	Mistral-24B
KMMLU (Son et al., 2025b)	0.710	0.744	0.589	0.501	0.551	0.572	0.565
HAE-RAE Bench (Son et al., 2024b)	0.695	0.742	0.658	0.542	0.577	0.609	0.606
CLLcK (Kim et al., 2024a)	0.815	0.836	0.670	0.620	0.712	0.675	0.713
KorMedMCQA (Kweon et al., 2024)	0.434	0.357	0.514	0.429	0.483	0.456	0.478
KBL (Kim et al., 2024b)	0.552	0.464	0.351	0.389	0.382	0.436	0.505

In Tables 12–13, we provide the per-benchmark accuracies for the datasets used in the test sets for Figures 6–8 in § 4.1.1. Among the included benchmarks, we report results for the top five benchmarks by coverage, as benchmarks for which only one or two samples do not provide fair accuracy comparisons.

G.2 CUSTOMIZED BENCHHUB

We provide three additional examples of real-world use cases of BENCHHUB:

- (c) **Legal chatbot servicing in Korea and the US:** To select a foundation model for a legal chatbot, we select English and Korean datasets whose fine-grained subject is law. The final accuracy is computed as an average of the English and Korean datasets, ensuring that the model holds legal knowledge in both countries.
- (d) **Docent agent for Korean traditional arts:** To identify the best-performing model with expertise in Korean traditional arts, we select Korean datasets within BENCHHUB whose fine-grained subjects are labeled as architecture, sculpture, and painting. To ensure balanced representation across individual subjects, the questions are drawn using a stratified sampling strategy at a subject level.
- (e) **Counseling agent servicing in Korea:** To evaluate counseling agent in Korean, we select Korean datasets comprising:
 1. psychology-related samples (*i.e.*, fine-grained category is psychology),
 2. samples aware to Korean social interactions (*i.e.*, coarse-grained category is social intelligence),
 3. samples relevant to common counseling topics (*i.e.*, fine-grained categories are work life, daily life, and family).

The final accuracy is computed as a weighted average of these subsets, with weights of 0.5, 0.3, and 0.2, respectively.

Table 14 presents the top-5 model rankings across these scenarios. The fluctuations in model rankings among the three scenarios also underscore the practical need for tailored evaluations using BENCHHUB.

Table 14: Top 5 LLMs evaluated by customized BENCHHUB across three scenarios

Rank	(c) Legal chatbot	(d) Docent for Korean art	(e) Counseling agent
1	Qwen3-32B	Qwen2.5-72B-Instruct	Qwen2.5-72B-Instruct
2	gemma-3-1b-it	gemma-3-27b-it	Qwen3-8B
3	Qwen3-8B	Llama-3.3-70B-Instruct	gemma-3-27b-it
4	Qwen3-1.7B	Qwen3-32B	DeepSeek-R1-Distill-Qwen-32B
5	Mistral-Small-24B-Instruct-2501	Mistral-Small-24B-Instruct-2501	Mistral-Small-24B-Instruct-2501

G.3 SIMULATION ON DATA CONTAMINATION

To assess the robustness of BENCHHUB under potential data contamination, we conduct a controlled simulation using OLMoE Muennighoff et al. (2025), a fully open-source model with publicly documented training data. We construct two variants of the model: (1) OLMoE-base, fine-tuned on 2k samples from the MATH Hendrycks et al. (2021b) training set, and (2) OLMoE-contaminated, fine-tuned on an equally sized subset of the MATH test set to emulate direct contamination. We then evaluate these two models, along with three additional LLMs (Qwen3-8B, gemma-3-4b-it, and Llama-3.1-8B-Instruct), on both the original MATH test set and BENCHHUB customized for math evaluation, which aggregates nine math-related benchmarks.

Table 15: Model ranking and accuracy across BENCHHUB customized for math evaluation and MATH (Hendrycks et al., 2021b) for simulation study on data contamination

Rank	BENCHHUB (Ours)	MATH (Hendrycks et al., 2021b)
1	Qwen/Qwen3-8B (0.35)	OLMoE-Contaminated(0.84)
2	gemma-3-4b-it(0.30)	Qwen/Qwen3-8B(0.31)
3	Llama-3.1-8B-Instruct (0.26)	gemma-3-4b-it (0.28)
4	OLMoE-Contaminated (0.22)	Llama-3.1-8B-Instruct (0.13)
5	OLMoE-base (0.16)	OLMoE-base (0.10)

Table 15 details the model accuracy and rankings on MATH and BENCHHUB customized for math evaluation. The results reveal stark differences in contamination sensitivity between single-benchmark and multi-benchmark evaluations. OLMoE-contaminated achieves extremely high accuracy (0.84) on the MATH test set, while OLMoE-base performs poorly (0.10), indicating that a single benchmark can be highly vulnerable to contamination. In contrast, BENCHHUB preserves the ranking of all five models and yields only a modest gap between the contaminated and uncontaminated OLMoE variants. This empirical evidence suggests that aggregating multiple benchmarks, as done in BENCHHUB, provides a more contamination-robust evaluation signal than relying on a single dataset.

H ADDITIONAL RELATED WORK

Table 16: Comparison to existing evaluation platforms

Method	Customization	Categorization			Dynamic Scalability	Evaluation Target
		Fine-Grained	Sample-Level	Automated		
FLASK Ye et al. (2024)	✗	✓	✓	✗	✗	LLM as a Generator
HELM Liang et al. (2023)	△	✗	✗	✗	✗	LLM as a Generator
RAVEL Huang et al. (2024)	✗	✓	✓	✗	✗	Interpretability Methods
LLMBar Zeng et al. (2024)	✗	✗	✗	✗	✗	LLM as a Evaluator
Arena Hard Li et al. (2024b)	✗	✗	✗	✗	✓	LLM as a Generator
BenchHub (Ours)	✓	✓	✓	✓	✓	LLM as a Generator

We summarize the key differences between our platform and existing evaluation platforms in Table 16.

HELM. HELM Liang et al. (2023) aims to support user-specific evaluation via scenario definitions, but its customization is limited to 16 fixed scenarios derived from single datasets (e.g., MATH). This fixed structure cannot capture the multi-domain, multi-criteria nature of real-world user intent.

2214 In contrast, BenchHub automatically constructs a benchmark suite tailored to arbitrary user intent,
2215 drawing from diverse datasets and categories (see Section 4.2).

2216 **FLASK.** While FLASK Ye et al. (2024) proposes a skill/domain taxonomy, its 38 categories mix
2217 skills and domains and cover only a subset of BenchHub’s taxonomy. Moreover, FLASK relies on
2218 manual annotation and is static, whereas BenchHub provides fully automated categorization across 64
2219 subjects, 3 skills, 2 cultural attributes, and 3 dataset-level attributes, enabling scalable customization.

2220 **Arena Hard.** Arena Hard Li et al. (2024b) introduces evaluation data derived from real-world
2221 prompts, but its domain coverage is limited to the topics present in its original source. While Arena
2222 Hard uses GPT-4 as a strong baseline for pairwise comparison, we demonstrate that no single LLM
2223 consistently dominates others; rankings fluctuate depending on the benchmark composition (see
2224 Section 4). Additionally, its evaluation relies heavily on proprietary LLM judges, which are known to
2225 exhibit preference bias. BenchHub instead focuses on unifying existing benchmarks according to
2226 user intent while maintaining model-agnostic evaluation.

2227 **RAVEL and LLMBAR.** We note that RAVEL Huang et al. (2024) and LLMBAR Zeng et al. (2024)
2228 target fundamentally different evaluation goals—interpretability methods and evaluator models,
2229 respectively—and therefore lie outside the scope of our benchmark-merging framework.

2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

I FULL EXPERIMENTAL RESULTS IN ACCURACY

See Table 17-18 for the scores (accuracies) of the models across subject types.

Table 17: Results of all models across fine-grained categories (English)

Subject	gpt-4.1	claude-3.7-sonnet	gemini-2.0	gemma-3-27b	DeepSeek-R1-32B	Llama-3.3-70B	Mistral-24B
Tech							
Urban eng.	0.882	0.765	0.824	0.625	0.765	0.588	0.882
Nuclear eng.	1.000	0.750	0.500	0.500	0.500	1.000	1.000
Marin eng.	1.000	0.667	1.000	0.500	1.000	1.000	1.000
Biomedical eng.	0.963	0.828	0.716	0.563	0.743	0.779	0.794
Mechanics	0.943	0.829	0.829	0.559	0.706	0.647	0.941
Materials eng.	0.987	0.920	0.760	0.595	0.811	0.784	0.932
IT	0.904	0.735	0.783	0.598	0.690	0.724	0.782
Environmental eng.	0.957	0.739	0.855	0.652	0.797	0.754	0.928
Energy eng.	0.953	0.802	0.791	0.628	0.826	0.767	0.872
Electrical eng.	0.877	0.816	0.825	0.609	0.722	0.704	0.800
Programming	1.000	0.913	0.826	0.667	0.611	0.556	0.722
Civil eng.	1.000	0.769	0.923	0.750	0.750	0.750	1.000
Chemical eng.	0.714	0.571	0.571	0.429	0.714	0.714	0.571
AI	0.931	0.984	0.817	0.474	0.420	0.355	0.330
Agricultural eng.	1.000	0.867	0.800	0.705	0.864	0.795	0.932
Aerospace eng.	1.000	0.833	1.000	1.000	0.833	0.833	1.000
Science							
Statistics	0.879	0.803	0.803	0.452	0.563	0.600	0.622
Physics	0.892	0.800	0.842	0.549	0.689	0.705	0.713
Mathematics	0.918	0.956	0.872	0.756	0.717	0.587	0.711
Life science	0.965	0.798	0.781	0.565	0.809	0.678	0.904
Geology	0.990	0.816	0.776	0.688	0.792	0.656	0.885
Earth science	0.979	0.798	0.840	0.692	0.788	0.779	0.942
Chemistry	0.863	0.814	0.762	0.510	0.650	0.697	0.720
Biology	0.959	0.730	0.818	0.533	0.767	0.769	0.835
Atmospheric science	0.990	0.753	0.753	0.739	0.783	0.641	0.935
Astronomy	0.965	0.843	0.843	0.704	0.835	0.809	0.852
HASS							
Welfare	0.896	0.722	0.729	0.576	0.654	0.737	0.797
Trade	0.944	0.807	0.800	0.494	0.811	0.767	0.856
Cognitive studies	0.620	0.524	0.481	0.500	0.580	0.662	0.629
Religion	0.912	0.877	0.895	0.724	0.914	0.860	0.948
Politics	0.909	0.759	0.693	0.635	0.767	0.767	0.872
Philosophy	0.875	0.664	0.632	0.455	0.711	0.623	0.651
Media	0.857	0.864	0.759	0.667	0.889	0.778	0.722
Literature	0.950	0.850	0.850	0.684	0.950	0.750	0.950
Law	0.750	0.596	0.610	0.294	0.540	0.518	0.679
Language	0.736	0.548	0.518	0.420	0.526	0.519	0.504
History	0.911	0.864	0.578	0.463	0.786	0.857	0.881
Geography	0.911	0.804	0.804	0.628	0.773	0.886	0.818
Education	0.957	0.793	0.793	0.580	0.795	0.652	0.848
Economics	0.893	0.809	0.695	0.574	0.597	0.713	0.752
Administration	0.899	0.797	0.732	0.551	0.819	0.819	0.841
Social Intelligence							
Value	0.699	0.890	0.788	0.653	0.599	0.857	0.619
Norms	0.816	0.658	0.605	0.516	0.613	0.581	0.710
Commonsense	0.837	0.765	0.749	0.871	0.877	0.856	0.837
Bias	0.000	1.000	0.333	0.349	0.333	0.324	0.288
Culture							
Work life	0.778	0.667	0.704	0.600	0.720	0.700	0.720
Tradition	0.833	0.881	0.950	0.618	0.806	0.800	0.784
Housing	1.000	1.000	0.750	1.000	1.000	0.750	0.750
Food	0.534	0.479	0.479	0.360	0.553	0.675	0.456
Family	0.913	0.739	0.609	0.591	0.659	0.705	0.818
Daily life	0.600	0.521	0.475	0.355	0.590	0.676	0.532
Clothing	1.000	1.000	1.000	1.000	1.000	1.000	1.000
Holiday	1.000	1.000	1.000	1.000	1.000	1.000	1.000
Arts % Sports							
Sports	0.781	0.578	0.453	0.714	0.929	0.786	0.857
Sculpture	1.000	1.000	1.000	0.500	1.000	0.500	1.000
Photography	1.000	0.600	0.800	0.400	0.400	0.800	0.800
Performing	0.846	0.846	0.769	0.673	0.654	0.808	0.846
Painting	1.000	0.600	0.900	0.600	0.900	0.700	1.000
Music	1.000	1.000	0.800	0.900	0.900	0.900	0.800
Festivals	0.500	1.000	1.000	1.000	0.500	1.000	0.500
Fashion	1.000	0.800	1.000	0.800	0.800	0.600	0.600
Architecture	1.000	0.857	0.714	0.429	1.000	0.571	1.000

Table 18: Results of all models across fine-grained categories (Korean)

Subject	gpt-4.1	claude-3.7-sonnet	gemini-2.0	gemma-3-27b	DeepSeek-R1-32B	Llama-3.3-70B	Mistral-24B
Tech							
Urban eng.	0.552	0.634	0.559	0.504	0.507	0.543	0.468
Nuclear eng.	0.676	0.647	0.618	0.676	0.559	0.588	0.588
Marine eng.	0.688	0.826	0.625	0.569	0.521	0.611	0.569
Biomedical eng.	0.838	0.805	0.409	0.727	0.507	0.767	0.713
Mechanics	0.661	0.709	0.563	0.537	0.495	0.487	0.420
Materials eng.	0.720	0.820	0.560	0.608	0.510	0.619	0.608
IT	0.854	0.877	0.667	0.727	0.756	0.803	0.742
Environmental eng.	0.591	0.649	0.480	0.456	0.427	0.462	0.368
Energy eng.	0.587	0.674	0.551	0.507	0.457	0.457	0.399
Electrical eng.	0.688	0.778	0.646	0.549	0.535	0.549	0.500
Programming	0.667	0.722	0.667	0.667	0.667	0.667	0.833
Civil eng.	0.517	0.669	0.530	0.503	0.391	0.497	0.430
Chemical eng.	0.711	0.809	0.641	0.596	0.539	0.574	0.560
AI	0.861	0.829	0.676	0.694	0.618	0.657	0.703
Agricultural eng.	0.605	0.605	0.539	0.464	0.386	0.506	0.428
Aerospace eng.	0.757	0.786	0.579	0.621	0.564	0.629	0.579
Science							
Statistics	0.813	0.813	0.571	0.571	0.582	0.549	0.615
Physics	0.826	0.870	0.644	0.626	0.595	0.603	0.542
Mathematics	0.842	0.889	0.848	0.385	0.487	0.359	0.359
Life science	0.783	0.783	0.635	0.635	0.609	0.739	0.635
Geology	0.755	0.765	0.627	0.608	0.422	0.618	0.510
Earth science	0.701	0.769	0.627	0.604	0.552	0.575	0.575
Chemistry	0.760	0.829	0.643	0.574	0.612	0.643	0.512
Biology	0.852	0.875	0.586	0.766	0.664	0.742	0.711
Atmospheric science	0.719	0.688	0.625	0.531	0.531	0.656	0.563
Astronomy	1.000	1.000	1.000	0.900	1.000	1.000	0.800
HASS							
Welfare	0.783	0.745	0.516	0.755	0.742	0.724	0.705
Trade	0.856	0.767	0.658	0.752	0.752	0.766	0.731
Religion	0.846	0.860	0.714	0.805	0.706	0.812	0.856
Psychology	1.000	1.000	1.000	1.000	1.000	1.000	0.000
Politics	0.806	0.858	0.714	0.717	0.634	0.667	0.703
Philosophy	0.843	0.897	0.715	0.791	0.718	0.757	0.757
Media	0.942	0.928	0.897	0.877	0.755	0.876	0.877
Literature	0.836	0.914	0.760	0.700	0.739	0.798	0.800
Law	0.604	0.555	0.463	0.510	0.416	0.544	0.530
Language	0.807	0.906	0.763	0.648	0.685	0.750	0.705
History	0.775	0.794	0.691	0.622	0.526	0.603	0.570
Geography	0.711	0.778	0.698	0.594	0.522	0.631	0.597
Education	0.732	0.816	0.586	0.701	0.603	0.755	0.660
Economics	0.814	0.820	0.606	0.704	0.701	0.692	0.656
Administration	0.731	0.766	0.598	0.691	0.635	0.711	0.675
Social Intelligence							
Value	0.848	0.879	0.697	0.818	0.818	0.788	0.758
Norms	0.884	0.881	0.881	0.881	0.810	0.721	0.762
Commonsense	0.835	0.873	0.822	0.718	0.757	0.748	0.767
Bias	0.993	0.966	0.951	1.000	1.000	0.846	1.000
Culture							
Work life	0.926	0.926	0.826	0.921	0.768	0.921	0.921
Tradition	0.962	0.960	0.858	0.917	0.819	0.900	0.911
Leisure	1.000	1.000	1.000	0.500	0.500	1.000	0.500
Housing	0.824	0.824	0.647	0.735	0.676	0.676	0.676
Food	0.850	0.923	0.769	0.744	0.684	0.789	0.821
Family	0.826	0.792	0.696	0.652	0.818	0.864	0.800
Daily life	0.837	0.837	0.823	0.751	0.682	0.738	0.764
Clothing	0.793	0.793	0.690	0.621	0.655	0.759	0.655
Holiday	0.643	0.602	0.602	0.620	0.616	0.674	0.654
Arts & Sports							
Sports	0.960	0.960	0.818	0.960	0.917	0.913	0.864
Sculpture	0.923	0.833	0.833	1.000	0.727	0.917	0.833
Photography	0.800	0.855	0.655	0.768	0.600	0.667	0.655
Performing	0.950	0.950	0.911	0.930	0.752	0.884	0.918
Painting	0.931	0.932	0.833	0.896	0.794	0.837	0.918
Music	0.912	0.971	0.758	0.909	0.667	0.879	0.909
Festivals	0.941	1.000	1.000	0.941	0.882	0.813	0.941
Fashion	0.626	0.626	0.524	0.565	0.490	0.571	0.456
Architecture	0.745	0.778	0.641	0.711	0.658	0.664	0.618

2376 Tables 19 and 20 details the accuracy of 14 open LLMs across four different sampling strategies in
 2377 English and Korean, respectively.

2378
 2379 Table 21 details the accuracy of 14 open LLMs evaluated by the customized BENCHHUB in five
 2380 different scenarios.

2381 Table 19: Evaluation results of 14 open LLMs in English across four different sampling strategies

2383	Model	Random	Stratified	Chatbot Arena	MixEval
2384	Qwen2.5-72B-Instruct	0.874	0.962	0.888	0.870
2385	Qwen3-1.7B	0.702	0.733	0.696	0.677
2386	Qwen3-14B	0.743	0.778	0.733	0.729
2387	Qwen3-32B	0.696	0.707	0.690	0.686
2388	Qwen3-4B	0.704	0.713	0.712	0.689
2389	Qwen3-8B	0.732	0.749	0.737	0.723
2390	DeepSeek-R1-Distill-Qwen-14B	0.729	0.763	0.726	0.707
2391	DeepSeek-R1-Distill-Qwen-32B	0.746	0.799	0.755	0.743
2392	gemma-3-1b-it	0.688	0.694	0.680	0.661
2393	gemma-3-27b-it	0.810	0.852	0.816	0.798
2394	gemma-3-4b-it	0.717	0.748	0.721	0.704
2395	Llama-3.1-8B-instruct	0.734	0.779	0.747	0.730
2396	Llama-3.3-70B-Instruct	0.784	0.833	0.788	0.779
2397	Mistral-Small-24B-Instruct-2501	0.817	0.845	0.811	0.789

2398 Table 20: Evaluation results of 14 open LLMs in Korean across four different sampling strategies

2400	Model	Random	Stratified	Chatbot Arena	MixEval
2401	Qwen2.5-72B-Instruct	0.360	0.376	0.363	0.371
2402	Qwen3-1.7B	0.624	0.647	0.646	0.630
2403	Qwen3-14B	0.697	0.708	0.723	0.692
2404	Qwen3-32B	0.597	0.613	0.547	0.606
2405	Qwen3-4B	0.671	0.674	0.683	0.666
2406	Qwen3-8B	0.605	0.562	0.615	0.518
2407	DeepSeek-R1-Distill-Qwen-14B	0.507	0.613	0.617	0.609
2408	DeepSeek-R1-Distill-Qwen-32B	0.531	0.533	0.510	0.541
2409	gemma-3-1b-it	0.661	0.666	0.665	0.649
2410	gemma-3-27b-it	0.453	0.474	0.469	0.468
2411	gemma-3-4b-it	0.613	0.635	0.638	0.625
2412	Llama-3.1-8B-instruct	0.612	0.623	0.618	0.623
2413	Llama-3.3-70B-Instruct	0.370	0.444	0.383	0.406
2414	Mistral-Small-24B-Instruct-2501	0.466	0.492	0.478	0.486

2415 Table 21: Evaluation results of 14 open LLMs using customized BENCHHUB across five use cases

2416	Model	(a)	(b)	(c)	(d)	(e)
2417	Qwen2.5-72B-Instruct	0.604	0.657	0.658	0.595	0.670
2418	Qwen3-1.7B	0.711	0.477	0.703	0.383	0.624
2419	Qwen3-4B	0.667	0.420	0.556	0.300	0.599
2420	Qwen3-8B	0.629	0.568	0.718	0.430	0.665
2421	Qwen3-14B	0.642	0.429	0.531	0.316	0.499
2422	Qwen3-32B	0.798	0.523	0.663	0.529	0.648
2423	DeepSeek-R1-Distill-Qwen-14B	0.657	0.554	0.653	0.479	0.647
2424	DeepSeek-R1-Distill-Qwen-32B	0.626	0.609	0.654	0.488	0.660
2425	Llama-3.1-8B-Instruct	0.650	0.581	0.602	0.393	0.627
2426	Llama-3.3-70B-Instruct	0.651	0.612	0.637	0.562	0.659
2427	Mistral-Small-24B-Instruct-2501	0.619	0.632	0.661	0.523	0.660
2428	gemma-3-1b-it	0.762	0.465	0.704	0.364	0.551
2429	gemma-3-4b-it	0.632	0.529	0.641	0.391	0.632
	gemma-3-27b-it	0.611	0.614	0.651	0.582	0.664