

P-MMEVAL: A Parallel Multilingual Multitask Benchmark for Consistent Evaluation of LLMs

Anonymous ACL submission

Abstract

Recent advancements in large language models (LLMs) showcase varied multilingual capabilities across tasks like translation, code generation, and reasoning. Previous assessments often limited their scope to fundamental natural language processing (NLP) or isolated capability-specific tasks. To alleviate this drawback, we aim to present a comprehensive multilingual multitask benchmark. First, we present a pipeline for selecting available and reasonable benchmarks from massive ones, addressing the oversight in previous work regarding the utility of these benchmarks, i.e., their ability to differentiate between models being evaluated. Leveraging this pipeline, we introduce P-MMEVAL, a large-scale benchmark covering effective fundamental and capability-specialized datasets. Furthermore, P-MMEVAL delivers consistent language coverage across various datasets and provides parallel samples. Finally, we conduct extensive experiments on representative multilingual model series to compare performances across models, analyze dataset effectiveness, examine prompt impacts on model performances, and explore the relationship between multilingual performances and factors such as tasks, model sizes, and languages, offering valuable guidance for future research ¹.

1 Introduction

In recent years, large language models (LLMs, Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023; Bai et al., 2022, 2023) have raised significant interest in the artificial intelligence (AI) community. As most LLMs are English-centric, when we focus on the performances of a specific LLM, it generally refers to the evaluation results on English benchmarks. For example, early research focuses on reporting evaluation results on fundamental natural language processing (NLP)

benchmarks. i.e, how accurately the LLM understands and generates text, including TRIVIAQA (Joshi et al., 2017a), WINOGRANDE (Sakaguchi et al., 2020), and HELLASWAG (Zellers et al., 2019). Nowadays, researchers are more interested in capability-specialized benchmarks, i.e., how well LLM performs on a group of specific task-solving problems, including GSM8K (Cobbe et al., 2021) for mathematical reasoning, MMLU (Hendrycks et al., 2021a) for knowledge acquisition, and HUMANEVAL (Chen et al., 2021) for code generation. However, there is currently little work on systematically evaluating the multilingual capabilities of LLMs. When developing and iterating LLMs, giving accurate and parallel evaluation results is crucial for identifying their multilingual capabilities and quantifying their performances.

Building a benchmark with both inclusive task coverage and strong linguistic parallelism is difficult. Measuring the multilingual abilities of a specific LLM, or comparing the quality of generated multilingual responses from one LLM to another, remains a big challenge in developing multilingual LLMs. Early work focuses on an isolated evaluation pipeline for a specific task, or to be more concrete, a specific perspective of LLM abilities: MHELLASWAG (Dac Lai et al., 2023) aims at collecting the multilingual understanding abilities, XLSUM (Hasan et al., 2021) mainly focus on evaluating the quality of generated multilingual text, HUMANEVAL-XL (Peng et al., 2024) is used for quantify how well-executed the generated code segments are, and MGSM (Shi et al., 2023) is made for testifying the performance on arithmetic reasoning. In modern research, for delivering simpler aggregation and comprehensive evaluation when judging model abilities, researchers collect several popular isolated benchmark tasks and propose a united, large-scale multilingual benchmark system like XTREME (Hu et al., 2020), XTREME-R (Ruder et al., 2021), XGLUE (Liang et al., 2020),

¹We will publish all the related resources in the future.

MEGA (Ahuja et al., 2023), and BUFFET (Asai et al., 2024) for multi-task assessments. However, these large-scale benchmarks 1) are tailored predominantly to fundamental NLP tasks and 2) inconsistently cover multiple languages across their selected datasets.

In this paper, our goal is to present a pipeline to develop a comprehensive multilingual multitask benchmark. To this end, we first select representative and challenging datasets from fundamental NLP tasks to reduce redundant testing and enhance the efficiency of evaluation. The second phase of our endeavor involves a meticulous curation of the most intensely studied capability-specialized tasks in contemporary research including code generation, knowledge comprehension, mathematical reasoning, logical reasoning, and instruction following. Finally, we construct a collection of datasets P-MMEVAL, consisting of three fundamental NLP datasets and five advanced capability-specialized datasets. To maintain language coverage among all selected datasets, we unify 10 languages considering the cost and computational limitations via expert translation review to construct the missing multilingual portions.

To summarize, our contributions are as follows:

- We present a pipeline for selecting available and reasonable benchmarks to assess the multilingual abilities of LLMs. Innovatively, we employ a statistical analysis method to identify effective datasets from a collection of ones. Our method can enhance the objectivity and scientific rigor of the selection process.
- We develop a multilingual multi-task benchmark P-MMEVAL that includes both fundamental and capability-specialized tasks, which ensures consistent language coverage across various datasets and provides parallel samples across different languages. This benchmark facilitates a thorough assessment of multilingual capabilities and enables unprecedented fairness and consistency in evaluating cross-lingual transfer capabilities.
- Our experiments offer a comprehensive analysis of the multilingual capabilities of various LLMs, showcasing performance across different prompts, models, languages, and tasks. Importantly, we analyze the utility of each dataset within P-MMEVAL in distinguishing model performance, thus identifying specific

benchmarks that differentiate model performance across model series and sizes.

2 Related Work

Isolated Fundamental NLP Benchmarks Although diverse multilingual evaluation benchmarks have been established, they focused on basic language understanding and generation capabilities of models. Notable work includes XNLI (Conneau et al., 2018) for natural language inference, XCOPA (Ponti et al., 2020), MHELLASWAG (Dac Lai et al., 2023), and XWINOGRAD (Tikhonov and Ryabinin, 2021) for commonsense reasoning, PAWS-X (Yang et al., 2019) for paraphrase identification, XL-WiC (Raganato et al., 2020) for word sense disambiguation, as well as the span extraction QA datasets including XQUAD (Artetxe et al., 2020), MLQA (Lewis et al., 2020), and TYDIQA-GOLDP (Joshi et al., 2017b). Additional examples include XLSUM (Hasan et al., 2021) for text summarization and FLORES-200 (Costa-jussà et al., 2022) for machine translation. Each of those benchmarks is typically designed for a specific task, solely focusing on one aspect of the model’s capabilities.

Unified Fundamental NLP Benchmarks There are also large-scale benchmarks that unify diverse existing datasets, aiming at offering a comprehensive evaluation of the model’s abilities from various perspectives. For instance, XTREME (Hu et al., 2020) comprises four tasks related to natural language understanding (NLU). Its refined version, XTREME-R (Ruder et al., 2021), optimizes the specific datasets tailored for each task category within XTREME. The XGLUE (Liang et al., 2020), MEGA (Ahuja et al., 2023), and BUFFET (Asai et al., 2024) benchmarks integrate various datasets for both understanding and generation tasks.

Capability-specialized Multilingual Benchmarks The advanced task-solving capabilities of LLMs have garnered significant attention from the research community. The six capabilities that receive the most emphasis are mathematical reasoning (Cobbe et al., 2021; Hendrycks et al., 2021b), logical reasoning (Liu et al., 2020), instruction following (Li et al., 2023), knowledge comprehension (Hendrycks et al., 2021a), code generation (Chen et al., 2021), and conversational abilities (Bai et al., 2024). Typical multilingual benchmarks include MGSM (Shi et al., 2023) for

Source	Task	Benchmarks	# Examples	Test sets	Metric	
Existing	Generation	FLORES-200 (Costa-jussà et al., 2022)	1012 × 10	Annotation	BLEU	
	Understanding	XNLI (Conneau et al., 2018) MHELLASWAG (Dac Lai et al., 2023)	120 × 10 (3) 120 × 10 (3)	Translation Translation	Acc Acc	
	Extension	Code generation	HUMANEVAL-XL (Peng et al., 2024)	80 × 10 (3) × 12	Translation	Pass@1
		Mathematical reasoning	MGSM (Shi et al., 2023)	250 × 10 (3)	Translation	Acc
		Logic reasoning	MLOGIQA (Liu et al., 2020)	80 × 10 (8)	Translation	Acc
		Knowledge	MMMLU (Hendrycks et al., 2021a)	400 × 10 (2)	Translation	Acc
		Instruction following	MIFEVAL (Zhou et al., 2023)	96 × 10 (9)	Translation	Acc

Table 1: An overview of the P-MMEVAL benchmark. In total, P-MMEVAL takes seven multilingual tasks into consideration, which is built on eight benchmarks. “# Examples” denotes “the number of examples per language” × “the number of involved languages” × “the number of programming languages” (special for HUMANEVAL-XL), and the numbers of extended languages are in parentheses. “Test sets” section describes the nature of the test sets (whether they are translations of English data or independently annotated).

180 mathematical reasoning, the OpenAI multilingual
 181 version of MMLU (MMMLU)² for knowledge
 182 comprehension, and HUMANEVAL-XL (Chen
 183 et al., 2021) for code generation.

184 All the benchmarks mentioned above focus ei-
 185 ther exclusively on fundamental NLP capabilities
 186 or on advanced application abilities. Additionally,
 187 there is inconsistent multilingual coverage across
 188 various datasets within a single multi-task bench-
 189 mark. The proposed benchmark P-MMEVAL in-
 190 tegrates three fundamental NLP datasets and five
 191 capability-specialized datasets, providing consis-
 192 tent language coverage across all selected datasets.

193 3 Datasets Selection Pipeline

194 Over time, the evaluation tasks for language mod-
 195 els encompass a wide variety, with each category
 196 amassing substantial multilingual datasets. These
 197 datasets are primarily categorized into two main
 198 types: generation and understanding. Each task is
 199 further divided into various subcategories, most of
 200 which consist of multiple datasets. Therefore, se-
 201 lecting effective ones is crucial, as it can reduce re-
 202 dundant testing and improve evaluation efficiency.

203 We suggest that benchmarks incapable of main-
 204 taining statistically significant performance differ-
 205 ences between models with notable capability dis-
 206 parities should be considered ineffective. Such
 207 benchmarks are often too simplistic, allowing even
 208 smaller models to achieve comparable performance.
 209 Consequently, these benchmarks fall short of pro-
 210 viding reliable and meaningful evaluations, as their
 211 difficulty levels are disproportionately lower than

212 the capabilities of current mainstream LLMs. Fur-
 213 thermore, in model optimization, if a benchmark is
 214 unable to provide significant differences in evalua-
 215 tion, it becomes difficult to effectively indicate the
 216 direction for improvement, making it challenging
 217 to adjust and enhance the model’s performance in
 218 a targeted manner. Therefore, selecting a bench-
 219 mark that can clearly distinguish between model
 220 performance is crucial for guiding the optimiza-
 221 tion process. Thus, we utilize paired-sample T-test
 222 (Field, 2005) to optimize the selection process by
 223 filtering out datasets that can effectively distin-
 224 guish the performances of LLMs with substantial size
 225 gaps. Empirically, we hypothesize that substantial
 226 capability differences exist between models with
 227 large size gaps within the same families. How-
 228 ever, due to differences in pre-training strategies
 229 and model architectures, it remains challenging to
 230 ensure consistent capability differences for models
 231 of the same size from different families. Therefore,
 232 the same size models of different series are not
 233 considered in the dataset selection process.

234 Our selection pipeline can be described as fol-
 235 lows: Given the evaluation results of model A and
 236 model B on a multilingual dataset D , denoted as
 237 A_i and B_i respectively, where i represents the lan-
 238 guage index. Following this, we first collect two
 239 score arrays $[A_1, A_2, \dots, A_m]$ and $[B_1, B_2, \dots, B_m]$
 240 which represents the evaluation results of model
 241 A and model B on m different languages, respec-
 242 tively. Then, we use these two arrays to derive
 243 the significance value p after running a paired-T
 244 significance test. If p is less than a pre-defined sig-
 245 nificance level (e.g., 0.01), it can be concluded that
 246 there is a significant difference in the overall scores

²<https://huggingface.co/datasets/openai/MMMLU>

247 between model *A* and model *B*. By determining
248 whether multiple pairs of models have significantly
249 different scores on this dataset, the effectiveness
250 of the dataset in distinguishing the performance
251 among various models can be identified.

252 4 P-MMEval

253 We aim to build a comprehensive evaluation system
254 that unifies diverse NLP and capability-specialized
255 tasks, ensures consistent language coverage per
256 task, and offers parallel samples across languages
257 to facilitate consistent comparisons. The overview
258 of our proposed P-MMEVAL is shown in Table 1.

259 4.1 Design Principles

260 **Diversity in tasks** First, the two key fundamental
261 NLP tasks of generating and understanding are
262 covered. More critically, through in-depth analysis,
263 we identify and establish five kinds of core capabilities
264 of current LLMs, including code generation,
265 knowledge comprehension, mathematical reasoning,
266 logical reasoning, and instruction following.

267 **Diversity in languages** To ensure that our benchmark
268 can also help testify the cross-lingual transferability
269 of LLMs, we unify 10 different languages spanning 7 language families, including
270 English (*en*), Chinese (*zh*), Arabic (*ar*), Spanish
271 (*es*), Japanese (*ja*), Korean (*ko*), Thai (*th*), French
272 (*fr*), Portuguese (*pt*), and Vietnamese (*vi*).

274 4.2 Fundamental NLP Dataset Curation

275 In light of the diversity of fundamental NLP
276 datasets, we meticulously select 11 datasets widely
277 employed in research (Ahuja et al., 2023; Asai
278 et al., 2024; Liang et al., 2020), spanning across
279 the two major categories of understanding and
280 generation. This curation aims to thoroughly appraise
281 the models’ foundational capabilities. Below, we
282 briefly summarize these two categories of tasks.

283 4.2.1 Tasks

284 **Natural Language Understanding (NLU)** Here,
285 we have five sub-tasks: i) The natural language inference
286 (NLI) dataset XNLI (Conneau et al., 2018),
287 ii) Three commonsense reasoning datasets encompass
288 XCOPA (Ponti et al., 2020), MHELLASWAG,
289 and XWINograd (Tikhonov and Ryabinin, 2021).
290 iii) The paraphrase identification dataset PAWS-X
291 (Yang et al., 2019). iv) The word sense disambiguation
292 dataset XL-WIC (Raganato et al., 2020).
293 v) Three span-prediction datasets, i.e., XQuAD

(Artetxe et al., 2020), MLQA (Lewis et al., 2020),
and TYDIQA-GOLDP (Joshi et al., 2017b).

Natural Language Generation (NLG) This task comprises the XLSUM (Hasan et al., 2021) and FLORES-200 (Costa-jussà et al., 2022) datasets. XLSUM is a multilingual summarization dataset derived from news articles. FLORES-200 is a dataset for multilingual machine translation, covering 200 languages.

302 4.2.2 Dataset Selection Process

Settings: We utilize three pairs of models to help fundamental benchmark curation, including Qwen2.5-7B vs. Qwen2.5-72B (Yang et al., 2024), LLAMA3.1-8B vs. LLAMA3.1-70B (Dubey et al., 2024), and MISTRAL-NEMO-INSTRUCT-2407 (MISTRAL-NEMO) vs. MISTRAL-LARGE-INSTRUCT-2407 (MISTRAL-LARGE).³ For understanding tasks, we utilize a fundamental prompt design with English instructions (See “EN” format in Section 5.2). For generation tasks, we employ the native prompt with instructions in the target language (See “Native” format in Section 5.2), as the “EN” prompt can cause the model to generate responses in English for non-English data. Then, we count the number of occurrences of each language in all benchmarks. For each benchmark, aside from English, we select four extra languages that are both supported in that benchmark and deserve the highest occurrences in all benchmarks. To expedite result verification, we gather a maximum of 250 instances per language across all tasks, ensuring an efficient yet comprehensive evaluation process.

Results: Table 2 presents the paired-sample T-test results, identifying significant differences in pairwise model performances on each dataset. The *p*-value threshold is set at 0.01. The dataset will be retained if all three selected model pairs show significant performance differences. Following this criterion, XNLI, MHELLASWAG, and FLORES-200 are retained for further processing.

335 4.3 Capability-specialized Dataset Curation

Besides the fundamental NLP tasks mentioned above, we also select one dataset for each of the five capability-specialized tasks.⁴ In detail, the

³<https://huggingface.co/mistralai/Mistral-Nemo-Instruct-2407> and <https://huggingface.co/mistralai/Mistral-Large-Instruct-2407>.

⁴For each specialized capability, we generally do not have enough choices (mostly only one benchmark is available).

Dataset	Available	Model series		
		QWEN	LLAMA	MISTRAL
<i>Understanding</i>				
XNLI	✓	0.0055	0.0009	0.0005
MHELLASWAG	✓	0.0028	0.0078	0.0039
PAWS-X	✗	0.5794	0.0170	0.0008
XL-WiC	✗	0.1734	0.0078	0.0058
XCOPA	✗	0.0070	0.0110	0.0014
XWINograd	✗	0.0224	0.0002	0.0014
XQuAD	✗	0.0283	0.0066	0.0117
TYDiQA-GOLDP	✗	0.2494	0.0375	0.0001
MLQA	✗	0.0011	0.0710	0.0064
<i>Generation</i>				
FLORES-200	✓	0.0010	0.0031	0.0007
XLSUM	✗	0.4835	0.7518	0.1500

Table 2: Results on significance test among three pairs of models: QWEN2.5-7B/72B (QWEN), LLAMA3.1-8B/70B (LLAMA), and MISTRAL-NEMO/LARGE (MISTRAL). For the understanding task and the generation task, we finally select XNLI and MHELLASWAG, and FLORES-200, respectively, as their significance level values are all lower than 0.01.

involved specialized capabilities in P-MMEVAL are:

- **Code generation** We utilize HUMANEVAL-XL (Peng et al., 2024) dataset, which establishes connections between 23 natural languages (NLs) and 12 programming languages (PLs).
- **Mathematical reasoning** We use the MGSM (Shi et al., 2023) dataset, a multilingual version translated from the monolingual GSM8K dataset consisting of math word problems.
- **Logical reasoning** We keep the original *en* and *zh* examples from origin LOGIQA (Liu et al., 2020) dataset.
- **Knowledge acquisition** We sample a subset of MMMLU comprising 200 “hard” samples and 200 “easy” samples. The performance of six diverse models (QWEN2.5-7B, QWEN2.5-72B, LLAMA3.1-8B, LLAMA3.1-70B, MISTRAL-NEMO, and MISTRAL-LARGE) is utilized as a proxy for selecting “hard” and “easy” samples. Concretely, we compile an “easy” subset comprising 6,335 instances where all models excel, and a “hard” subset consisting of 663 instances that challenge every model. Subsequently, guided by annotations from MMLU-REDUX (Gema et al., 2024), we

refine these subsets by discarding 798 erroneous instances from the “easy” pool and 160 from the “hard” pool. Finally, we systematically sample 200 instances from each of the pruned pools, thus creating our finalized “easy” and “hard” evaluation sets.

- **Instruction following** We employ the English IFEVAL (Liu et al., 2020) dataset, which consists examples following pre-defined 25 types of “verifiable instruction”.

4.4 Expansion of the Selected Datasets

To maintain consistency across all languages, we extend the support of some benchmark datasets on the missing languages by collecting human-annotated translation results. The number of expanded languages and samples for each dataset is listed in the “#Example” column of Table 1. More details of sampling are provided in Appendix ??.

We initially generate translated examples using the advanced GPT-4O⁵ model. Subsequently, a professional translation team conducts an exhaustive review of the machine translation outputs, correcting any errors, localizing vocabulary, and removing instances that do not translate well across languages. This meticulous process ensures both high translation quality and cultural adaptability.

The modification rate by post-review is detailed in Table 7. It is apparent that datasets contain translation errors to varying extents, with error rates peaking at 82.50%. This underscores the limitations of using raw machine-generated translations for dataset extension, highlighting the critical need for human review to maintain translation fidelity. Notably, among the most frequent errors are mis-translations of proper nouns and inconsistencies in terminology usage, followed by omissions. These trends indicate that the model currently struggles with specific domain terminology and maintaining contextual coherence.

4.5 Significance Detection for the Expanded Datasets

In the previous sections, we utilize a meticulous dataset selection framework to pinpoint three essential datasets: XNLI, MHellaSwag, and Flores-200. Furthermore, we broaden the range of languages covered in each dataset, with the additional languages specified in parentheses in the “#Example” column of Table 1. We follow

⁵gpt-4o-2024-05-13

Model	Understanding		Code generation	Mathematical reasoning	Logic reasoning	Knowledge	Instruction following	Generation	AVG_S	AVG_U
	XNLI	MHELLASWAG								
<i>Open-source models (<7B)</i>										
LLAMA3.2-1B	31.67	24.49	37.71	12.08	27.12	27.80	35.42	29.30	28.03	28.08
LLAMA3.2-3B	30.67	23.74	37.42	11.64	25.62	26.85	34.90	36.85	27.29	27.21
QWEN2.5-0.5B	22.25	19.68	33.92	13.12	14.62	30.25	30.21	15.95	24.42	20.97
QWEN2.5-1.5B	46.58	36.35	48.59	35.20	35.12	42.02	44.37	21.37	41.06	41.47
QWEN2.5-3B	60.08	48.09	60.75	69.40	39.38	46.27	66.46	25.75	56.45	54.09
GEMMA2-2B	53.50	45.31	51.54	44.52	34.88	40.85	56.67	24.00	45.69	49.41
<i>Open-source models (7-14B)</i>										
LLAMA3.1-8B	52.84	49.11	69.96	67.24	39.88	43.80	59.27	16.59	56.03	50.98
QWEN2.5-7B	67.17	62.92	71.88	81.08	45.88	49.83	77.71	32.76	65.28	65.05
GEMMA2-9B	57.92	65.62	69.96	81.28	41.50	49.23	79.17	36.48	64.23	61.77
MISTRAL-NEMO	54.25	55.73	57.38	76.52	41.75	44.88	60.00	33.65	56.11	54.99
QWEN2.5-14B	67.50	70.10	72.83	88.68	53.50	51.52	79.48	31.31	69.20	68.80
<i>Open-source models (14-50B)</i>										
QWEN2.5-32B	68.33	76.38	75.88	90.88	57.38	52.27	83.33	32.13	71.95	72.36
GEMMA2-27B	68.00	64.12	76.67	85.28	50.50	49.42	81.35	42.23	68.64	66.06
<i>Open-source models (>50B)</i>										
LLAMA3.1-70B	63.17	67.25	74.75	88.28	52.38	55.52	79.17	16.63	70.02	65.21
QWEN2.5-72B	71.42	75.95	76.00	91.00	58.38	52.67	87.60	41.55	73.13	73.69
MISTRAL-LARGE	69.58	69.04	77.17	90.48	53.50	51.85	83.23	43.40	71.25	69.31
<i>Closed-source models</i>										
GPT-4O	69.17	81.04	77.05	91.60	56.75	55.77	85.21	46.32	73.28	75.11
CLAUDE-3.5-SONNET	71.50	77.72	82.92	92.84	62.25	56.17	80.73	16.20	74.98	74.61

Table 3: Evaluation results of different models on P-MMEVAL. We gather those models by referring to their sizes. AVG_U and AVG_S represent the average score of the understanding and capability-specialized tasks, respectively. HUMAN EVAL-XL score presents the average score of three programming languages.

the setting in subsection 4.2.2 to utilize the paired-sample T-test to verify the effectiveness of the expanded datasets. Significance test results can be found in Appendix Table 5. The findings demonstrate that the expanded datasets continue to proficiently distinguish performance disparities among models with notable capability differences, thereby affirming the robustness and validity of our dataset selection framework.

4.6 Instruction selection

We utilize English instructions from OPENCOMPASS (Contributors, 2023) and LM-EVALUATION-HARNESS (Dac Lai et al., 2023). Among multiple instructions, we select a suitable one and make uniform modifications to ensure consistency across similar tasks. For zero-shot prompts, to increase the success rate of answer extraction, we add a constraint at the end of the instruction to some tasks, requiring the model to output the generated answers in a fixed format. In addition, we translate English instructions into multiple languages to construct native instructions.

5 Experiments

This section focuses on the following aspects: assessing the multilingual capabilities of different models; assessing the utility of each dataset within P-MMEVAL in distinguishing model performance;

examining the influence of various prompts on multilingual performance; and analyzing the correlation between models’ performance in English and non-English languages. All evaluation results are presented in Table 3.

5.1 Multilingual Models

We evaluate the performance of several representative instruction-tuned models – (i) closed-source models GPT-4O⁶ (OpenAI, 2023) and CLAUDE-3.5-SONNET⁷, (ii) open-source models including LLAMA3.1, LLAMA3.2 (Dubey et al., 2024), QWEN2.5 (Yang et al., 2024), MISTRAL-NEMO, MISTRAL-LARGE, and GEMMA2 series (Rivière et al., 2024).

5.2 Evaluation Settings

According to Zhao et al. (2021), the choice of prompts significantly impacts the evaluation results of LLMs and the model performance is sensitive to minor variations in prompting. In this study, we compare the evaluation results using the following prompts. EN: Instructions in English + input in the target language. Native: Instructions in the target language + input in the target language. EN-Few-Shot: Instructions in English + demonstrations in the target language + input in the target language.

⁶gpt-4o-2024-05-13

⁷claude-3-5-sonnet-20240620

467 For MGSM, we employ Chain of Thought (CoT)
468 ([Wei et al., 2022](#)) reasoning, which guides the
469 model to think step-by-step before providing a fi-
470 nal answer. For the other datasets, direct answer-
471 ing is utilized, which requests the model to pro-
472 duce answers directly. The inference methods for
473 these datasets align with the most commonly used
474 settings. Notably, for MMMLU, we choose the
475 prompt template following OpenAI simple-evals
476 repository.⁸ Specifically, CoT reasoning exhibits a
477 significantly higher answer extraction failure rate
478 compared to direct answering on small-sized LLMs
479 (i.e., the number of parameters is less than 7B),
480 leading to poor performance. Thus, we employ a
481 direct answering prompt for small-sized LLMs.⁹

482 For the few-shot demonstrations, we primarily
483 sample demonstrations from the validation set. For
484 the missing multilingual portions, we utilize GPT-
485 4O to translate these demonstrations from English
486 into the missing languages.

487 5.3 Main Results

488 Table 3 presents an overview of the evaluation re-
489 sults. Unless otherwise noted, the standard EN
490 prompt is applied to all datasets except FLORES-
491 200, HUMAN EVAL-XL, and MIFEVAL, where
492 the Native prompt is required. More information
493 about the prompting strategies including EN, Na-
494 tive, and En-Few-Shot is shown in Appendix C.
495 The evaluation result on HUMAN EVAL-XL is the
496 average score across three programming languages
497 including Python, JavaScript, and Java. See Ap-
498 pendix G for programming language evaluation
499 details. For the Flores-200 dataset, in addition to
500 reporting BLEU scores, we also provide COMET
501 scores measured by wmt22-comet-da ([Rei et al.](#))
502 (see Appendix, Table 8).

503 First, the multilingual capabilities of models be-
504 come stronger as the model sizes increase ([Kaplan](#)
505 et al., 2020). One exception is that when the size
506 of LLAMA3.2 increases from 1B to 3B, there is
507 a slight decline in performance. The main reason
508 for this is that LLAMA3.2-1B and LLAMA3.2-
509 3B exhibit poor instruction-following capabilities,
510 leading to a higher failure rate in answer extraction
511 and, consequently, fluctuations in the final score.
512 As the model size increases, the improvements in
513 various multilingual tasks show significant differ-
514 ences. Evaluation results on the understanding and

⁸<https://github.com/openai/simple-evals>

⁹The detailed evaluation prompts are illustrated in Appendix K.

515 capability-specialized tasks show significant im-
516 provement in understanding context, processing
517 semantic information, reasoning, and special abil-
518 ities, with increasing model sizes. For example,
519 for the QWEN2.5 series, the scores on the MGSM
520 dataset for the 0.5B and 72B models are 13.12 and
521 91.00, respectively. In contrast, the models’ per-
522 formance on generation tasks is relatively weaker
523 and shows slight improvement. Evaluations on
524 the FLORES-200 datasets indicate that, despite the
525 increase in model size, the generation capability
526 does not improve proportionally. This may reflect
527 the complexity of generating text that maintains
528 logical coherence and contextual relevance, where
529 increasing model sizes does not significantly en-
530 hance output quality.

531 In addition, QWEN2.5 demonstrates a strong
532 multilingual performance on understanding and
533 capability-specialized tasks, while GEMMA2 ex-
534 cels in generation tasks. CLAUDE-3.5-SONNET
535 performs poorly on FLORES-200 because it tends
536 to generate additional relevant statements in its
537 responses, potentially downgrading the BLEU
538 score. GPT-4O generally outperforms open-source
539 models. The performance gap between the best-
540 performing open-source model and GPT-4O is
541 within 3%.

542 6 Analyses

543 6.1 Analysis on Dataset Utility

544 The primary objective of this section is to assess
545 the utility of each dataset within P-MMEVAL in
546 distinguishing model performances. We divide
547 open-sourced models into categories by two as-
548 pects: model series and model sizes. Specifically,
549 we collect 5 categories of models from 5 model
550 series: QWEN2.5 (0.5B, 1.5B, 3B, 7B, 14B, 32B,
551 72B), LLAMA3.1 (8B, 70B), LLAMA3.2 (1B,
552 3B), GEMMA2 (2B, 9B, 27B), MISTRAL (Nemo,
553 Large).

554 And, we divide them into three categories based
555 on their sizes: Less than 7B (QWEN2.5-0.5B,
556 QWEN2.5-1.5B, QWEN2.5-3B, LLAMA3.2-1B,
557 LLAMA3.2-3B, GEMMA2-2B), Between 7B and
558 14B (QWEN2.5-7B, LLAMA3.1-8B, GEMMA2-
559 9B, MISTRAL-NEMO, QWEN2.5-14B), Larger
560 than 70B (LLAMA3.1-70B, QWEN2.5-72B,
561 MISTRAL-LARGE).

562 Table 4 shows the utility of each dataset in distin-
563 guishing the performances of paired models within
564 the same category. The detailed method for cal-

Dataset	MISTRAL	LLAMA3.2	LLAMA3.1	QWEN2.5	GEMMA2	>70B	7B-14B	<7B
FLORES-200	2/2	1/2	2/2	4/7	3/3	3/3	2/5	1/6
MHELLASWAG	2/2	1/2	2/2	6/7	2/3	2/3	5/5	5/6
XNLI	2/2	1/2	2/2	5/7	3/3	2/3	3/5	5/6
HUMANEVAL-XL (Python)	2/2	1/2	2/2	2/7	1/3	3/3	3/5	3/6
HUMANEVAL-XL (JavaScript)	2/2	1/2	2/2	5/7	3/3	2/3	5/5	5/6
HUMANEVAL-XL (Java)	2/2	1/2	2/2	4/7	3/3	2/3	3/5	3/6
MGSM	2/2	1/2	2/2	6/7	3/3	1/3	4/5	4/6
MLOGIQA	2/2	1/2	2/2	6/7	3/3	2/3	3/5	3/6
MIFEVAL	2/2	1/2	2/2	6/7	2/3	3/3	2/5	4/6

Table 4: All tested models are categorized into 8 categories based on model size or series. This table presents the utility of each dataset in distinguishing the performances of paired models within the same category. A value closer to 1 indicates higher utility for the dataset, with a value of 1 signifying that all models demonstrate distinguishable performances. Conversely, a numerator of 1 indicates that no models are distinguishable on that dataset. We set the threshold at 0.5, where each value is considered effective or ineffective in distinguishing the performances of models with the specified dataset.

culating the utility of each dataset is presented in Appendix I. A value closer to 1 indicates higher utility for the dataset, with a value of 1 signifying that all models within the same category demonstrate distinguishable performances. Conversely, a numerator of 1 indicates that no models are distinguishable on that dataset. We set the utility threshold at 0.5, where each value is considered effective or ineffective in distinguishing the performances of models with the specified dataset. Based on the results in Table 4, we can draw the following conclusions: 1) LLAMA3.2-1B and LLAMA3.2-3B show no significant performance differences across all datasets, indicating similar multilingual capabilities. The performance differentiation of small-size models below 7B is slightly worse. 2) Compared to JavaScript and Java, most models show poor performance differentiation in Python. According to the Appendix G, the average score of all the tested open-source models in Python is 90.46, significantly higher than the scores in the other two languages (48.95 and 46.66, respectively), indicating that all models have a strong knowledge grasp in Python. 3) All selected datasets can distinguish between models in the majority of categories, which verifies the effectiveness of all datasets included in P-MMEVAL.

6.2 Performances on English vs. Non-English Benchmarks

To preliminarily explore the relationship between non-English ability and English ability of the model, we evaluate the performance of the QWEN2.5 series models (7B, 14B, 32B, and 72B) on six datasets with parallel samples. For each dataset, we calculate the ratio of the average score achieved on the test sets in all nine non-English

languages to the score achieved on the test data in English. We do not consider models smaller than 7B, as these models are easily influenced by prompts, leading to performance fluctuations.

Appendix Figure 1 illustrates the trend of the ratio of non-English performance to English performance as model sizes increase. On five datasets, the model’s non-English performance appears limited by its English performance. However, on the three programming languages (Python, JavaScript, Java) of HUMANEVAL-XL dataset, the models achieve comparable performance in both English and non-English test sets. This means that code knowledge is less dependent on natural language. When the model size increases, we observe that: 1) As for instruction-following ability, the gap between non-English data and English data is narrowing. 2) The ratio of capability-specialized datasets outperforms those of fundamental understanding datasets.

7 Conclusion

In this paper, we first present a pipeline for benchmark selection, which guides the finding and selecting of effective benchmarks for quantifying the multilingual performances of LLMs. Then, we introduce a comprehensive multilingual multitask benchmark, P-MMEVAL, which covers both fundamental and capability-specialized tasks, ensuring consistent language coverage and providing parallel samples in multiple languages. Furthermore, we conduct extensive experiments on representative multilingual model series. These findings provide valuable guidance for future research, highlighting the importance of balanced and comprehensive training data, effective prompt engineering, and the need for targeted improvements in specific language capabilities.

637 Limitations

638 Through the above experiments and analyses, we
639 summarize the following limitations:

640 1) Language Coverage: While P-MMEval cur-
641 rently covers 10 languages from 7 language fami-
642 lies, there is a need to include more languages to
643 better represent global linguistic diversity. Future
644 work will focus on expanding the language cover-
645 age to ensure a more comprehensive evaluation of
646 multilingual LLMs.

647 2) Task Diversity: P-MMEval includes 7 rep-
648 resentative tasks, but the rapidly evolving field of
649 LLMs demands a broader range of tasks. Future
650 work will focus on expanding the benchmark to
651 cover more diverse and challenging tasks, provid-
652 ing a more thorough assessment of multilingual
653 LLMs.

654 Ethics Statement

655 All procedures performed in studies involving hu-
656 man participants were in accordance with the eth-
657 ical standards of the institutional and/or national
658 research committee and with the 1964 Helsinki
659 Declaration and its later amendments or compara-
660 ble ethical standards. This article does not contain
661 any studies with animals performed by any of the
662 authors. Informed consent was obtained from all
663 individual participants included in the study.

664 References

665 Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent
666 Ochieng, Krithika Ramesh, Prachi Jain, Akshay Ut-
667 tama Nambi, Tanuja Ganu, Sameer Segal, Mohamed
668 Ahmed, Kalika Bali, and Sunayana Sitaram. 2023.
669 MEGA: multilingual evaluation of generative AI. In
670 *Proceedings of the 2023 Conference on Empirical
671 Methods in Natural Language Processing, EMNLP
672 2023, Singapore, December 6-10, 2023*, pages 4232–
673 4267. Association for Computational Linguistics.

674 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.
675 2020. On the cross-lingual transferability of mono-
676 lingual representations. In *Proceedings of the 58th
677 Annual Meeting of the Association for Compu-
678 tational Linguistics, ACL 2020, Online, July 5-10, 2020*,
679 pages 4623–4637. Association for Computational
680 Linguistics.

681 Akari Asai, Sneha Kudugunta, Xinyan Yu, Terra
682 Blevins, Hila Gonen, Machel Reid, Yulia Tsvetkov,
683 Sebastian Ruder, and Hannaneh Hajishirzi. 2024.
684 BUFFET: benchmarking large language models for
685 few-shot cross-lingual transfer. In *Proceedings of
686 the 2024 Conference of the North American Chap-
687 ter of the Association for Computational Linguistics:*

688 *Human Language Technologies (Volume 1: Long
689 Papers), NAACL 2024, Mexico City, Mexico, June
690 16-21, 2024*, pages 1771–1800. Association for Com-
691 putational Linguistics.

692 Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jia-
693 heng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su,
694 Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024.
695 Mt-bench-101: A fine-grained benchmark for eval-
696 uating large language models in multi-turn dialogues.
697 In *Proceedings of the 62nd Annual Meeting of the
698 Association for Computational Linguistics (Volume
699 1: Long Papers), ACL 2024, Bangkok, Thailand, Au-
700 gust 11-16, 2024*, pages 7421–7454. Association for
701 Computational Linguistics.

702 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
703 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
704 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
705 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
706 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
707 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
708 Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-
709 guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang,
710 Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,
711 Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx-
712 uan Zhang, Yichang Zhang, Zhenru Zhang, Chang
713 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang
714 Zhu. 2023. Qwen technical report. *arXiv preprint
arXiv:abs/2309.16609*.

716 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
717 Amanda Askell, Jackson Kernion, Andy Jones, Anna
718 Chen, Anna Goldie, Azalia Mirhoseini, Cameron
719 McKinnon, Carol Chen, Catherine Olsson, Christo-
720 pher Olah, Danny Hernandez, Dawn Drain, Deep
721 Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez,
722 Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua
723 Landau, Kamal Ndousse, Kamile Lukosiute, Liane
724 Lovitt, Michael Sellitto, Nelson Elhage, Nicholas
725 Schiefer, Noemí Mercado, Nova DasSarma, Robert
726 Lasenby, Robin Larson, Sam Ringer, Scott John-
727 ston, Shauna Kravec, Sheer El Showk, Stanislav Fort,
728 Tamera Lanham, Timothy Telleen-Lawton, Tom Con-
729 erly, Tom Henighan, Tristan Hume, Samuel R. Bow-
730 man, Zac Hatfield-Dodds, Ben Mann, Dario Amodei,
731 Nicholas Joseph, Sam McCandlish, Tom Brown, and
732 Jared Kaplan. 2022. Constitutional AI: harmlessness
733 from AI feedback. *arXiv preprint arXiv:2212.08073*.

734 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
735 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
736 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
737 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
738 Gretchen Krueger, Tom Henighan, Rewon Child,
739 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
740 Clemens Winter, Christopher Hesse, Mark Chen, Eric
741 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
742 Jack Clark, Christopher Berner, Sam McCandlish,
743 Alec Radford, Ilya Sutskever, and Dario Amodei.
744 2020. Language models are few-shot learners. In *Ad-
745 vances in Neural Information Processing Systems 33:
746 Annual Conference on Neural Information Process-
747 ing Systems 2020, NeurIPS 2020, December 6-12,
748 2020, virtual*.

749	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Koplán, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebbgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code. <i>arXiv preprint arXiv:abs/2107.03374</i> .	808 arXiv e-prints, pages arXiv–2307.
750		809
751		
752		
753		
754		
755		
756		
757		
758		
759		
760		
761		
762		
763		
764		
765		
766		
767		
768		
769		
770	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. <i>arXiv preprint arXiv:2110.14168</i> .	810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842
771		
772		
773		
774		
775		
776	Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. XNLI: evaluating cross-lingual sentence representations. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018</i> , pages 2475–2485. Association for Computational Linguistics.	843 844
777		
778		
779		
780		
781		
782		
783		
784	OpenCompass Contributors. 2023. Opencompass: A universal evaluation platform for foundation models. https://github.com/open-compass/opencompass .	845 846 847 848 849 850 851 852 853
785		
786		
787		
788	Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Y. Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barraud, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Rogers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. <i>arXiv preprint arXiv:abs/2207.04672</i> .	854 855 856 857 858 859 860 861
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804	Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A Rossi, and Thien Huu Nguyen. 2023. Okapi: Instruction-tuned large language models in multiple languages	862 863 864 865 866 867
805		
806		
807		
808	with reinforcement learning from human feedback. <i>arXiv e-prints</i> , pages arXiv–2307.	
809		
810	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Greigerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaïdis, Damien Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vraneš, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasudevan Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. The llama 3 herd of models. <i>arXiv preprint arXiv:abs/2407.21783</i> .	
811		
812		
813		
814		
815		
816		
817		
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828		
829		
830		
831		
832		
833		
834		
835		
836		
837		
838		
839		
840		
841		
842		
843	Andy Field. 2005. Discovering statistics using ibm spss statistics. <i>Sage</i> .	
844		
845	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-ku Lo, Craig Stewart, George F. Foster, Alon Lavie, and Ondrej Bojar. 2021. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In <i>Proceedings of the Sixth Conference on Machine Translation, WMT@EMNLP 2021, Online Event, November 10-11, 2021</i> , pages 733–774. Association for Computational Linguistics.	
846		
847		
848		
849		
850		
851		
852		
853		
854	Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, Claire Barale, Robert McHardy, Joshua Harris, Jean Kadour, Emile van Krieken, and Pasquale Minervini. 2024. Are we done with mmlu? <i>arXiv preprint arXiv:abs/2406.04127</i> .	
855		
856		
857		
858		
859		
860		
861		
862	Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Samin Mubashir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. Xl-sum: Large-scale multilingual abstractive summarization for 44 languages. In <i>Findings of the Association for Computational Linguistics</i>	
863		
864		
865		
866		
867		

868	<i>tics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021, volume ACL/IJCNLP 2021 of Findings of ACL, pages 4693–4703. Association for Computational Linguistics.</i>	924
869		925
870		926
871		927
872	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding. In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.	928
873		929
874		930
875		931
876		932
877		933
878	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the MATH dataset. In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	934
879		935
880		936
881		
882		
883		
884		
885		
886	Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalization. <i>arXiv preprint arXiv:abs/2003.11080</i> .	937
887		938
888		939
889		940
890		941
891	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017a. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers</i> , pages 1601–1611. Association for Computational Linguistics.	942
892		
893		
894		
895		
896		
897		
898		
899	Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017b. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers</i> , pages 1601–1611. Association for Computational Linguistics.	943
900		944
901		
902		
903		
904		
905		
906		
907	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>CoRR</i> , abs/2001.08361.	945
908		946
909		947
910		948
911		949
912	Patrick S. H. Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. MLQA: evaluating cross-lingual extractive question answering. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 7315–7330. Association for Computational Linguistics.	950
913		951
914		952
915		
916		
917		
918		
919	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval .	953
920		954
921		955
922		956
923		957
924	Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Dixin Jiang, Guihong Cao, Xiaodong Fan, Ruofei Zhang, Rahul Agrawal, Edward Cui, Sining Wei, Taroon Bharti, Ying Qiao, Jiun-Hung Chen, Winnie Wu, Shuguang Liu, Fan Yang, Daniel Campos, Rangan Majumder, and Ming Zhou. 2020. XGLUE: A new benchmark dataset for cross-lingual pre-training, understanding and generation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 6008–6018. Association for Computational Linguistics.	958
925		959
926		960
927		
928		
929		
930		
931		
932		
933		
934		
935		
936		
937	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. In <i>Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020</i> , pages 3622–3628. ijcai.org.	961
938		962
939		963
940		964
941		965
942		966
943	OpenAI. 2023. GPT-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	967
944		968
945	Qiwei Peng, Yekun Chai, and Xuhong Li. 2024. Humaneval-xl: A multilingual code generation benchmark for cross-lingual natural language generalization. In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy</i> , pages 8383–8394. ELRA and ICCL.	969
946		970
947		971
948		972
949		973
950		974
951		975
952		976
953	Edoardo Maria Ponti, Goran Glavas, Olga Majewska, Qianchu Liu, Ivan Vulic, and Anna Korhonen. 2020. XCOPA: A multilingual dataset for causal commonsense reasoning. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 2362–2376. Association for Computational Linguistics.	977
954		978
955		979
956		980
957		981
958		
959		
960		
961	Alessandro Raganato, Tommaso Pasini, José Camacho-Collados, and Mohammad Taher Pilehvar. 2020. Xlwic: A multilingual benchmark for evaluating semantic contextualization. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 7193–7206. Association for Computational Linguistics.	982
962		983
963		984
964		985
965		986
966		987
967		988
968		989
969	Ricardo Rei, José G. C. de Souza, Duarte M. Alves, Chrysoula Zerva, Ana C. Farinha, Taisiya Glushkova, Alon Lavie, Luísa Coheur, and André F. T. Martins. COMET-22: unbabel-ist 2022 submission for the metrics shared task. In <i>Proceedings of the Seventh Conference on Machine Translation, WMT 2022, Abu Dhabi, United Arab Emirates (Hybrid)</i> , December 7-8, 2022, pages 578–585.	990
970		991
971		992
972		993
973		994
974		995
975		996
976		997
977	Morgane Rivière, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Huszenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos,	998
978		999
979		
980		
981		

982	Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Patterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozinska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucinska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju-yeong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjösund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, and Lilly McNealus. 2024. <i>Gemma 2: Improving open language models at a practical size.</i> <i>arXiv preprint arXiv:abs/2408.00118.</i>	1042
983		1043
984		1044
985		1045
986		1046
987		
988	<i>Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021, volume ACL/IJCNLP 2021 of Findings of ACL, pages 3534–3546. Association for Computational Linguistics.</i>	
989		
990	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenjin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madiam Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Miaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <i>Llama 2: Open foundation and fine-tuned chat models.</i> <i>arXiv preprint arXiv:2307.09288.</i>	1047
991		1048
992		1049
993		1050
994		1051
995		1052
996		1053
997		1054
998		1055
999		1056
1000		1057
1001		1058
1002		1059
1003		1060
1004		1061
1005		1062
1006		1063
1007		1064
1008		1065
1009		1066
1010		1067
1011		1068
1012		1069
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020		
1021	Sebastian Ruder, Noah Constant, Jan A. Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu, Junjie Hu, Dan Garrette, Graham Neubig, and Melvin Johnson. 2021. XTREME-R: towards more challenging and nuanced multilingual evaluation. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 10215–10245. Association for Computational Linguistics.	1070
1022		1071
1023		1072
1024		1073
1025		1074
1026		1075
1027		1076
1028		1077
1029		
1030		
1031	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 8732–8740. AAAI Press.	1078
1032		1079
1033		1080
1034		1081
1035		1082
1036		1083
1037		1084
1038		1085
1039	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multi-lingual chain-of-thought reasoners. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1086
1040		1087
1041		1088
1042		1089
1043		1090
1044		1091
1045		1092
1046		1093
1047		1094
1048		
1049		
1050		
1051		
1052		
1053		
1054		
1055		
1056		
1057		
1058		
1059		
1060		
1061		
1062		
1063		
1064		
1065		
1066		
1067		
1068		
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		
1080		
1081		
1082		
1083		
1084		
1085		
1086		
1087		
1088		
1089		
1090		
1091		
1092		
1093		
1094		
1095		
1096		
1097		
1098		
1099		
1100		
1101		

November 3-7, 2019, pages 3685–3690. Association
for Computational Linguistics.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
Farhadi, and Yejin Choi. 2019. Hellaswag: Can a
machine really finish your sentence? In *Proceedings
of the 57th Conference of the Association for Compu-
tational Linguistics, ACL 2019, Florence, Italy, July
28- August 2, 2019, Volume 1: Long Papers*, pages
4791–4800. Association for Computational Linguis-
tics.

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and
Sameer Singh. 2021. Calibrate before use: Improv-
ing few-shot performance of language models. In
*Proceedings of the 38th International Conference on
Machine Learning, ICML 2021, 18-24 July 2021, Vir-
tual Event*, volume 139 of *Proceedings of Machine
Learning Research*, pages 12697–12706. PMLR.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Sid-
dhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou,
and Le Hou. 2023. [Instruction-following evalua-
tion for large language models](#). *arXiv preprint
arXiv:abs/2311.07911*.

	QWEN	LLAMA	MISTRAL
XNLI	2.70E-04	1.80E-04	1.10E-06
MHELLASWAG	4.30E-06	3.50E-06	1.40E-05
FLORES-200	1.70E-04	1.60E-03	3.10E-04

Table 5: Results on significance test among three pairs
of models: QWEN2.5-7B/72B (QWEN), LLAMA3.1-
8B/70B (LLAMA), and MISTRAL-NEMO/LARGE
(MISTRAL).

A Significance Detection for the Selected Datasets

Table 5 presents the significance test results of the
three expanded datasets.

Dataset	Native	EN	EN-Few-shot
MMMLU	44.30	44.69	45.70
MLOGIQA	42.27	41.96	44.88
MGSM	62.13	64.17	63.28
MHELLASWAG	52.03	53.37	59.07
XNLI	54.49	55.31	64.08
FLORES-200	30.00	24.31	29.18

Table 6: Comparison on P-MMEVAL using three differ-
ent prompt settings.

B Sampling Process for Each Dataset in P-MMEVAL

Specifically, since FLORES-200 already includes
data for 10 languages, no additional translation

was required. We retain the complete test set for
evaluation.

For HUMAN-EVAL-XL and MGSM, which con-
tain 80 and 250 examples per language respectively,
we ensured comprehensive coverage by translating
the entire set for each language.

For single-task datasets XNLI, MHELLASWAG,
and MLOGIQA, with large available test data, we
follow established practices and select the first N
examples for translation. This approach aligns with
prior literature (Shi et al., 2023) and ensures consist-
ency while managing computational and resource
constraints.

For multi-task datasets such as MMMLU and
IFEVAL, we adopt different strategies. For
MMMLU, we sample a subset comprising 200
“hard” samples and 200 “easy” samples, by utiliz-
ing diverse model evaluation results as a proxy.
For IFEVAL, we select 10 examples per task type,
resulting in a total of 110 examples. During the
translation verification process, 14 examples were
removed due to quality issues, leaving a final set of
96 examples.

C The Impact of Different Prompts on Model Performance

We explore three different prompting strategies:
EN, Native, and En-Few-Shot. Table 6 illustrates
the average performance of all evaluated open-
source models on various datasets of P-MMEVAL.
Overall, the performance difference between the
EN prompt and the Native prompt is minimal, re-
maining within 2%, indicating no substantial per-
formance gap. However, in the case of the FLORES-
200, the EN prompt results in a marked decline in
performance compared to the Native prompt. We
observe that models always generate responses in
English when English instructions are used to de-
scribe the task for non-English data for generation
tasks. On various datasets, the few-shot prompt
leads to better model performance than the zero-
shot prompt, as models achieve a higher success
rate in extracting answers in the few-shot setting.

D Comparison of Non-English and English Performance

Fig. 1 illustrates the ratio of non-English per-
formance to English performance with increasing
model sizes of QWEN2.5.

Dataset	<i>zh</i>	<i>ar</i>	<i>es</i>	<i>ja</i>	<i>ko</i>	<i>th</i>	<i>fr</i>	<i>pt</i>	<i>vi</i>
XNLI	/	/	/	22.50	11.67	/	/	10.83	/
MHELLASWAG	/	/	/	82.50	77.50	26.67	/	/	/
HUMANEval-XL	/	/	/	42.50	23.75	31.25	/	/	/
MGSM	/	9.20	/	/	32.80	/	/	5.60	27.20
MLOGIQA	/	22.50	30.00	51.25	33.75	46.25	3.75	46.25	18.75
MMMLU	/	/	/	/	/	26.00	13.50	/	/
MIFEVAL	25.50	23.81	20.00	45.71	36.19	37.14	21.90	17.14	24.76

Table 7: The table presents the percentage of modifications made by professional translators to the machine translation results.

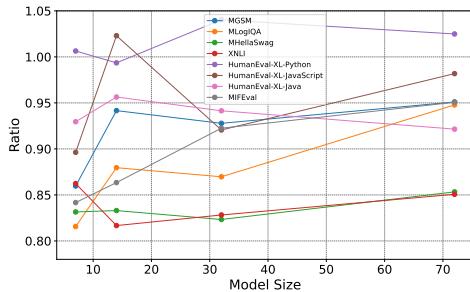


Figure 1: Illustration on the ratio of non-English performance to English performance with increasing model sizes of QWEN2.5.

E Expert Translation Review Results on Each Dataset

To supplement the missing multilingual portions in each dataset, a strategy that combines machine translation with professional human review is adopted. Table 7 shows the percentage of modifications made by professional translators to the machine translation results generated by GPT-4o. The main types of translation errors include omissions, incorrect translation order, and improper use of localized vocabulary.

F Evaluation of COMET Scores on the Flores-200 Dataset

In addition to the BLEU scores, we also provide COMET scores measured using the wmt22-comet-da model, shown in Table 8. For all tested models, the COMET scores are significantly higher than the BLEU scores, indicating that COMET is a more forgiving evaluation metric. Unlike BLEU, which requires strict literal matching, COMET focuses more on the semantics and fluency of the translation. This allows it to more comprehensively reflect translation quality and the models’ performance across different languages.

Model	COMET	BLEU
LLaMA3.2-1B	86.45	29.30
LLaMA3.2-3B	88.40	36.85
Qwen2.5-0.5B	88.83	15.95
Qwen2.5-1.5B	81.16	21.37
Qwen2.5-3B	80.58	25.75
Gemma2-2B	87.16	24.00
LLaMA3.1-8B	88.27	16.59
Qwen2.5-7B	87.75	32.76
Gemma2-9B	88.76	36.48
Mistral-Nemo	80.06	33.65
Qwen2.5-14B	85.17	31.31
Qwen2.5-32B	87.08	32.13
Gemma2-27B	87.62	42.23
LLaMA3.1-70B	87.26	16.63
Qwen2.5-72B	88.56	41.55
Mistral-Large	88.88	43.40

Table 8: The table displays the comparison between BLEU and COMET scores on the Flores-200 dataset.

G Evaluation Results on Three Programming Languages of HumanEval-XL

Table 9 shows the evaluation results of all tested models on three programming languages of HumanEval-XL. Model performance in Python greatly exceeds the performance in the other two programming languages. For instance, Gemma2-2B scores 98.13 in Python, compared to 29.25 in JavaScript and 27.25 in Java. Additionally, as the model size increases, there is a noticeable improvement in performance for both JavaScript and Java.

H Model performance on each language with Increasing Model Sizes

This section analyzes the trend of the performance of the model in each language with increasing model sizes. We only report the average performance on four capability-specialized datasets (HumanEval-XL, MGSM, MLogiQA, and MIFEVAL). In addition, we do not consider models

1179
1180

1181
1182
1183
1184
1185
1186
1187
1188
1189

1190
1191

1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202

1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222

	Python	JavaScript	Java
LLAMA3.2-1B	92.13	9.38	11.63
LLAMA3.2-3B	91.50	9.75	11.00
QWEN2.5-0.5B	78.38	14.25	9.13
QWEN2.5-1.5B	81.63	35.88	28.25
QWEN2.5-3B	84.00	53.75	44.50
GEMMA2-2B	98.13	29.25	27.25
LLAMA3.1-8B	96.38	46.88	66.63
QWEN2.5-7B	86.75	68.00	60.88
GEMMA2-9B	98.75	54.63	56.50
MISTRAL-NEMO	93.25	39.63	39.25
QWEN2.5-14B	84.50	72.75	61.25
QWEN2.5-32B	89.38	73.13	65.13
GEMMA2-27B	99.63	63.75	66.63
LLAMA3.1-70B	98.75	63.38	62.13
QWEN2.5-72B	85.63	75.00	67.38
MISTRAL-LARGE	88.63	73.88	69.00
GPT-4O	89.13	77.88	64.13
CLAUDE-3.5-SONNET	99.75	74.00	75.00

Table 9: The table presents the performance on three programming languages of HumanEval-XL.

smaller than 7B, as these models are easily influenced by prompts, leading to performance fluctuations. Model performance varies by language, with English demonstrating the strongest capabilities, while Thai and Japanese show the weakest.

I Dataset Utility

To quantify the utility of each dataset, we employ paired-sample T-tests for each pair of models within the same categories. Inspired by (Freitag et al., 2021), our main motivation is to try to divide models in the same category into several groups based on their pairwise significance gaps, where all model pairs in the same group do not have significant performance gaps, and performances of all model pairs from different groups are hard to be fully distinguished. Given the list of all models $\mathbf{m} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_m]$, we recurrently gather some of the models into the same group $\Omega_i = \{\mathbf{m}_{\pi_1}, \mathbf{m}_{\pi_2}, \dots, \mathbf{m}_{\pi_k}\}$, $\pi_j \in [1, 2, \dots, m]$ for $j \in [1, 2, \dots, k]$ at the i -th step, where: 1) for each model \mathbf{m}_{π_j} in Ω_i , it does not have a significant performance gap against any

model in Ω_i except itself:

$$f_1 = \begin{cases} \text{true if } \mathcal{T}(\mathbf{m}_{\pi_j}, \mathbf{m}_{\pi_p}) > \theta \text{ holds for any} \\ \quad p \in [1, 2, \dots, k], j \neq p; \\ \text{false otherwise;} \end{cases} \quad (1)$$

2) for each model in Ω_i , it has significant performance gaps against all the model not in Ω_i :

$$f_2 = \begin{cases} \text{true if } \mathcal{T}(\mathbf{m}_{\pi_j}, \mathbf{m}_p) < \theta \text{ holds for all} \\ \quad p \notin [\pi_1, \pi_2, \dots, \pi_k]; \\ \text{false otherwise;} \end{cases} \quad (2)$$

where $\mathcal{T}(\cdot, \cdot)$ returns the p -value of the performances between two given models, and θ represents the threshold for denoting significance level. The group Ω_i is fixed if f_1 and f_2 both hold true. Such a recurrent process continues till each model is gathered into one specific group.¹⁰

After gathering all models into several groups, we use the ratio of the number of such groups to the number of models to describe the utility of the specific dataset. A higher ratio means that we have more gathered groups, indicating that the benchmark is of high utility in distinguishing the performances of models. On the contrary, a lower ratio means that most of the models can be gathered into the same group, denoting that the benchmark may hardly tell which model performs better than any other model.

The algorithm for quantifying the utility of each benchmark dataset is presented in Algorithm 1.

J Significance Detection on Each Dataset

The section illustrates the significant difference between models' pairwise performance for all categories of models.

K The Prompt Utilized for Each Dataset

The section presents the inference prompt utilized for each dataset.

¹⁰See Algorithm 1 in Appendix I for more details.

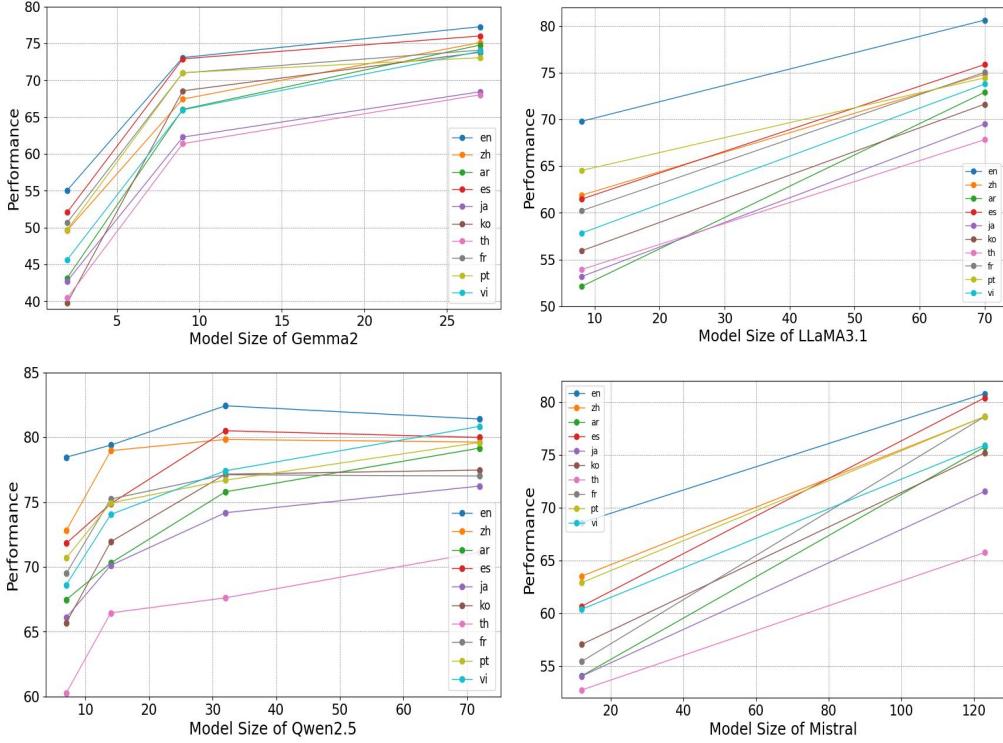


Figure 2: This figure illustrates the trend of the performance of the model in each language with increasing model sizes.

Algorithm 1 Algorithm for Quantifying the Utility of a Specific Benchmark Dataset

Input: Model ids $\mathbf{m} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_m]$, paired-sample T-test p -values among all pairs of models $p_{\mathbf{m}_i, \mathbf{m}_j} \in \mathbb{R}$ ($0 \leq i, j \leq m$, $p_{ij} = p_{ji}, i \neq j$), significance threshold $\theta \in \mathbb{R}$

Output: The number of sets $|\Omega|$, where Ω is a list of sets $\Omega = [\Omega_1, \Omega_2, \dots, \Omega_s]$, and each set contains several models $\Omega_i = \{\mathbf{m}_{\pi_1}, \mathbf{m}_{\pi_2}, \dots, \mathbf{m}_{\pi_k}\}$, $\Omega_i \neq \emptyset$, $|\Omega| = k \leq m$, $\pi_j \in [1, 2, \dots, m]$ for $j \in [1, 2, \dots, k]$

```

1:  $\Omega \leftarrow []$  ▷ Initialize with an empty list
2:  $\mathbf{z} = \mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_m$ 
3: while  $\mathbf{z} \neq \emptyset$  do
4:    $\mathbf{x} \leftarrow \{\mathbf{z}_1\}$  ▷ Initialize the current set with the first model id
5:    $\mathbf{y} \leftarrow \mathbf{z} - \mathbf{x}$ 
6:   while  $\mathbf{y} \neq \emptyset$  do
7:     Initialize  $\Gamma$  as a matrix full of  $\phi$ 
8:     for  $c \in \mathbf{x}$  do
9:       for  $d \in \mathbf{y}$  do
10:        if  $p_{c,d} < \theta$  then
11:           $\Gamma[c, d] \leftarrow \text{true}$ 
12:           $\Gamma[d, c] \leftarrow \text{true}$  ▷ The gap is significant
13:        else
14:           $\Gamma[c, d] \leftarrow \text{false}$ 
15:           $\Gamma[d, c] \leftarrow \text{false}$  ▷ The gap is not significant
16:        if  $\Gamma[c, d] = \text{false}$  for any  $c \in \mathbf{x}, d \in \mathbf{y}$  then ▷ Some paired models do not have significant performance gaps
17:          for  $d \in \mathbf{y}$  do
18:            if  $\Gamma[c, d] = \text{false}$  for any  $c \in \mathbf{x}$  then
19:               $\mathbf{x} \leftarrow \mathbf{x} + \{d\}$  ▷ Moving model  $d$  into the same group
20:               $\mathbf{y} \leftarrow \mathbf{y} - \{d\}$ 
21:            else ▷ Each model from  $\mathbf{x}$  has significant gap against each model from  $\mathbf{y}$ 
22:               $\Omega \leftarrow \Omega + [\mathbf{x}]$  ▷ Appending the new group  $\mathbf{x}$  into  $\Omega$ 
23:               $\mathbf{z} \leftarrow \mathbf{z} - \mathbf{x}$  ▷ Removing the processed model ids from  $\mathbf{z}$ 
24: return  $|\Omega|$  ▷ Return the number of groups

```

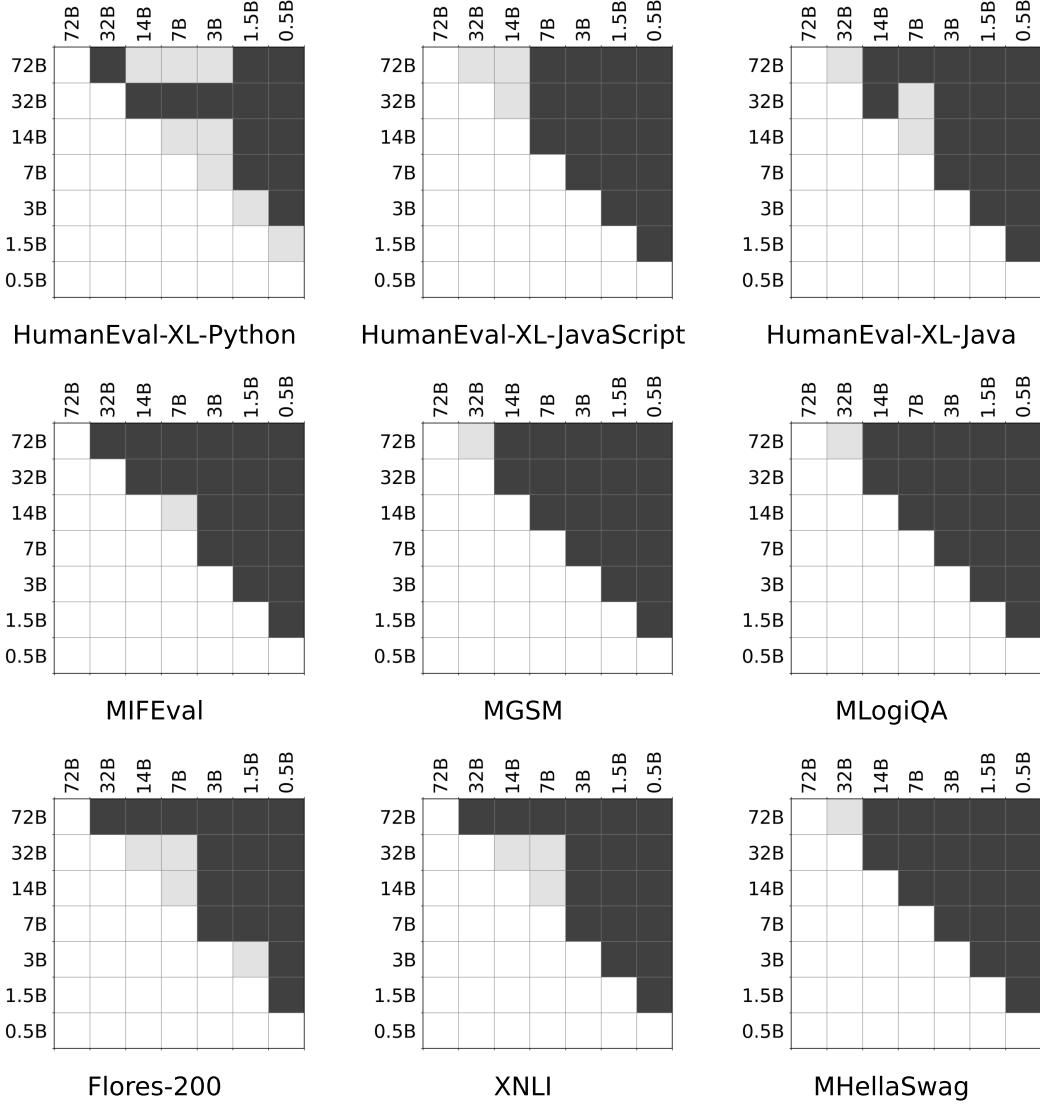


Figure 3: This figure illustrates the significant difference in pairwise performance among QWEN2.5 series models. Black blocks indicate that the p -values of paired t-tests between the corresponding models (vertical and horizontal) are less than 0.01, while gray blocks indicate p -values greater than 0.01.

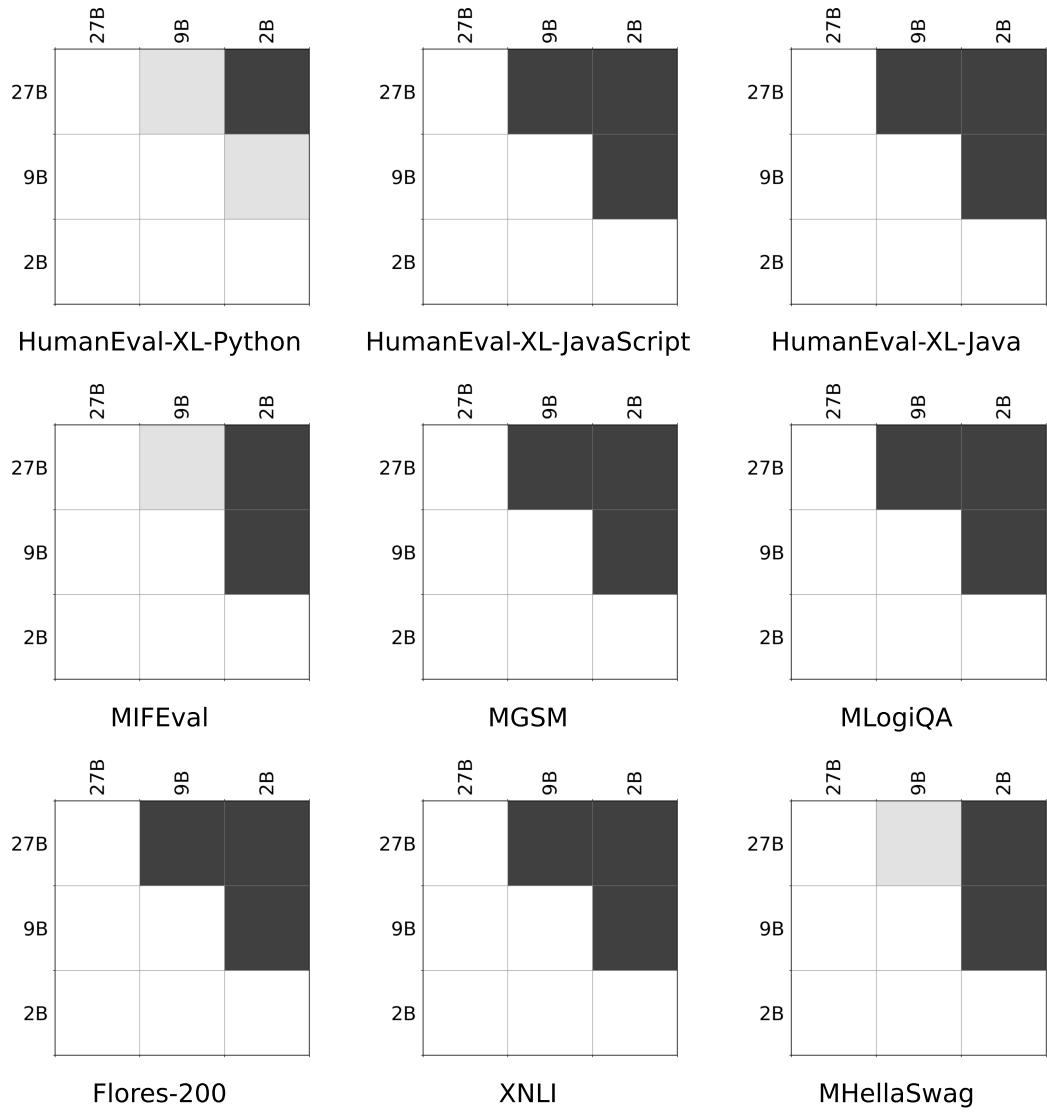


Figure 4: This figure illustrates the significant difference in pairwise performance among GEMMA2 series models.

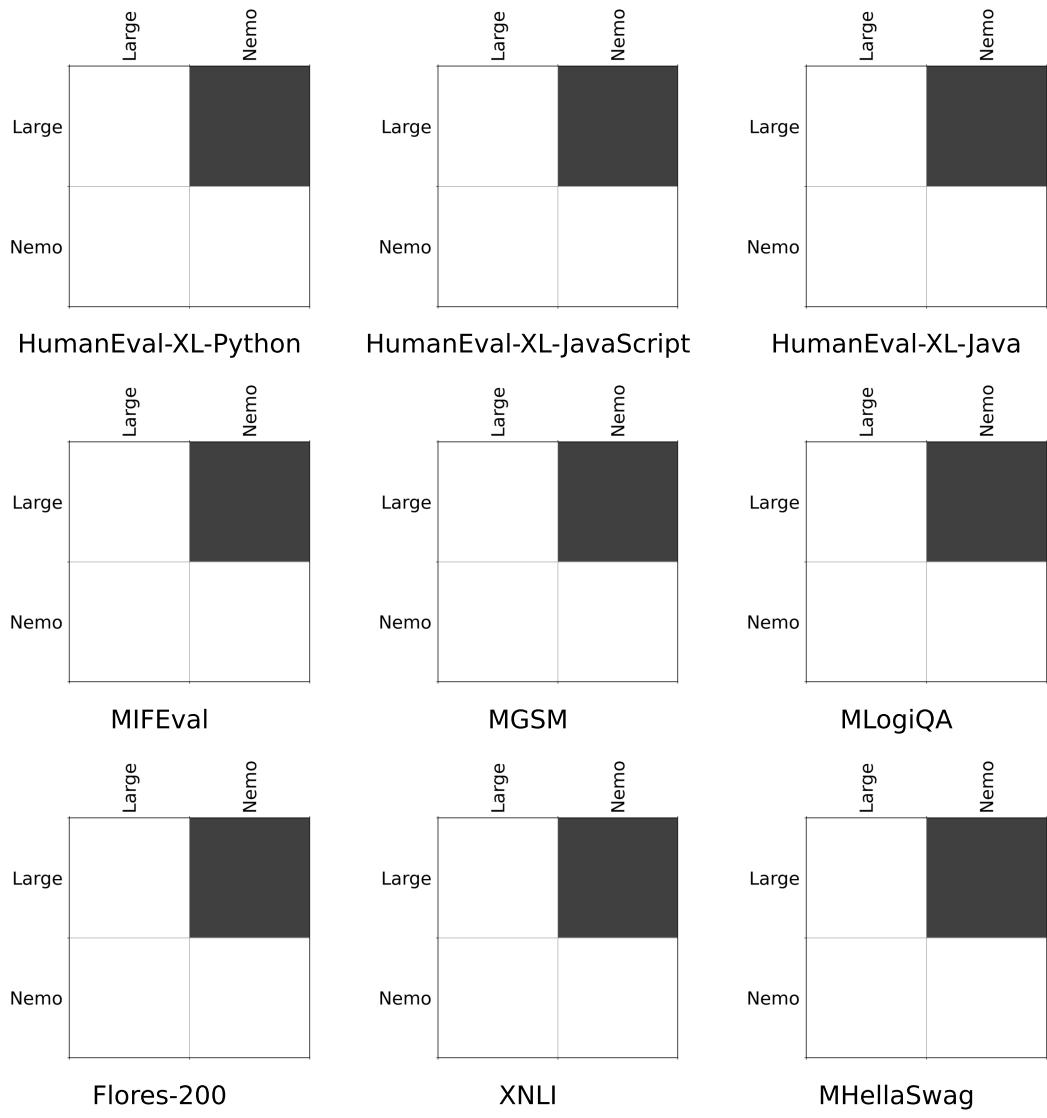


Figure 5: This figure illustrates the significant difference in pairwise performance among MISTRAL series models.

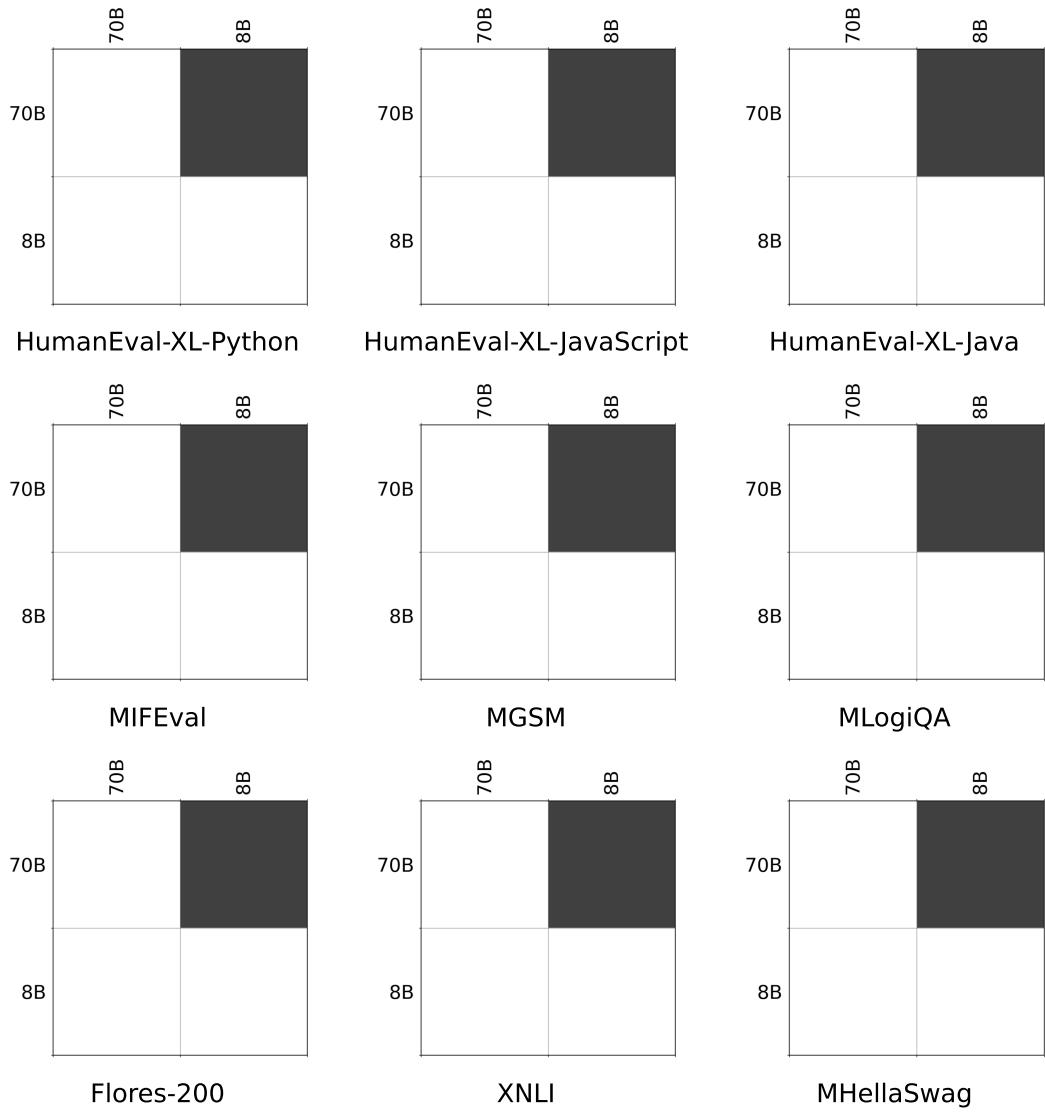


Figure 6: This figure illustrates the significant difference in pairwise performance among LLAMA3.1 series models.

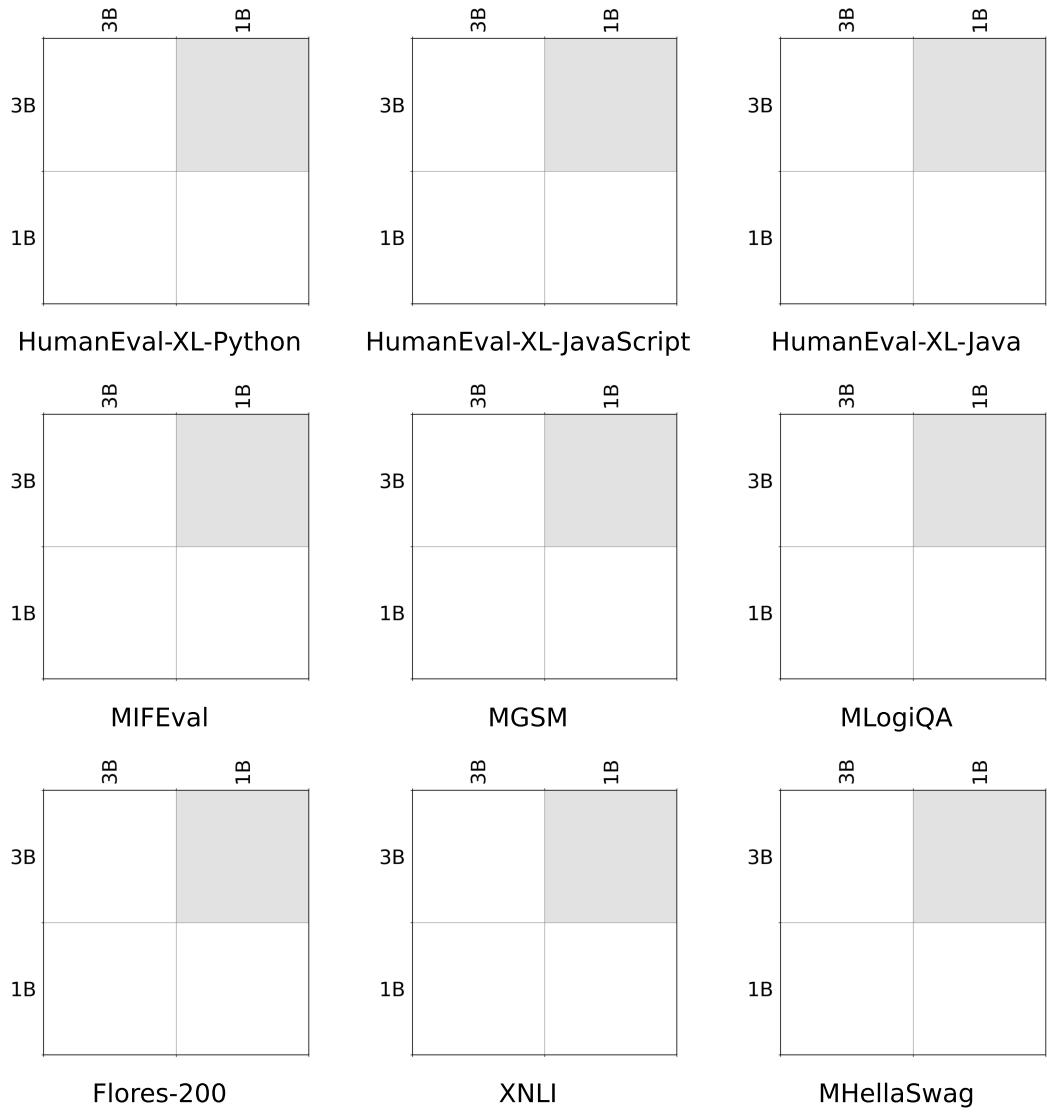


Figure 7: This figure illustrates the significant difference in pairwise performance among LLAMA3.2 series models.

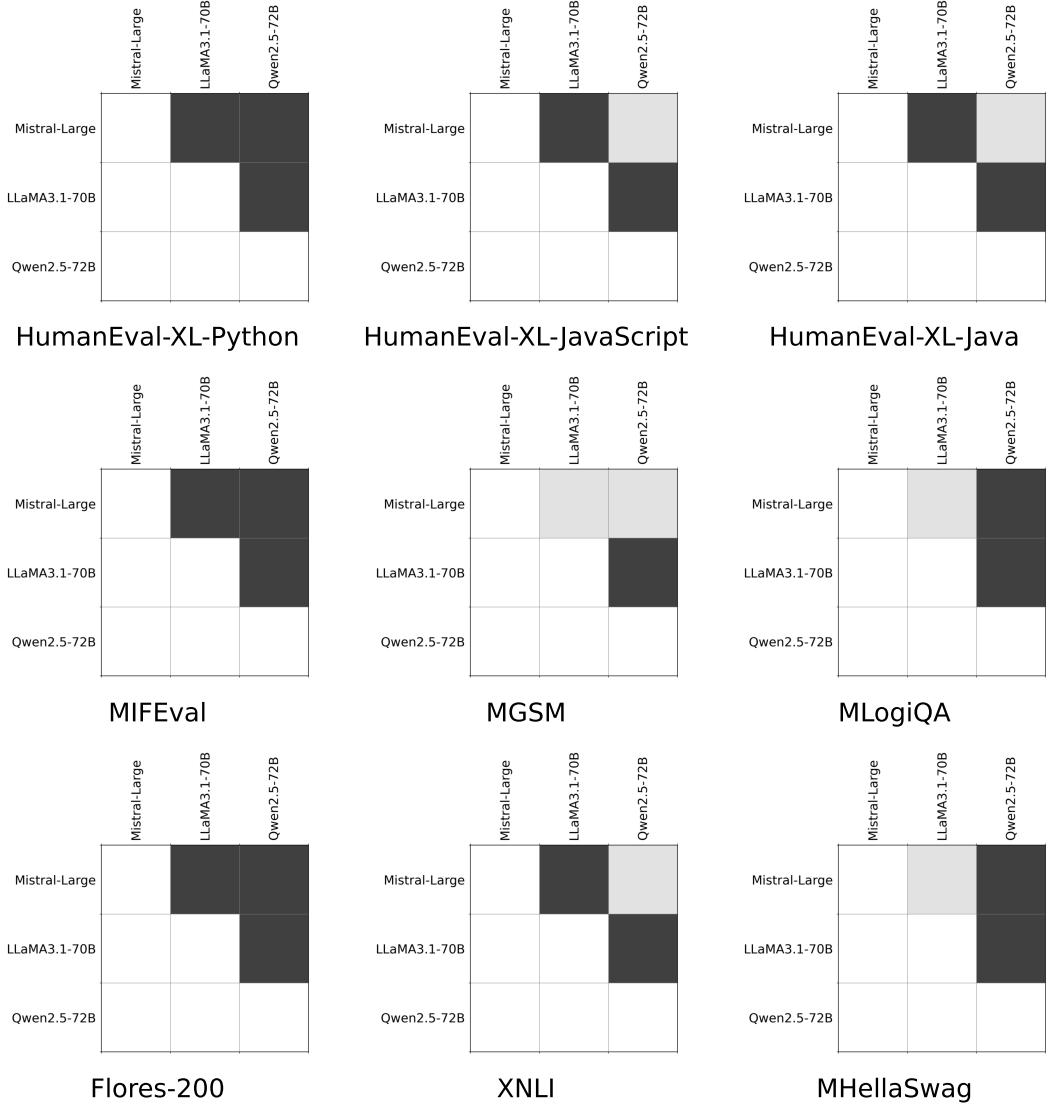


Figure 8: This figure illustrates the significant difference in pairwise performance among models with more than 70 billion parameters.

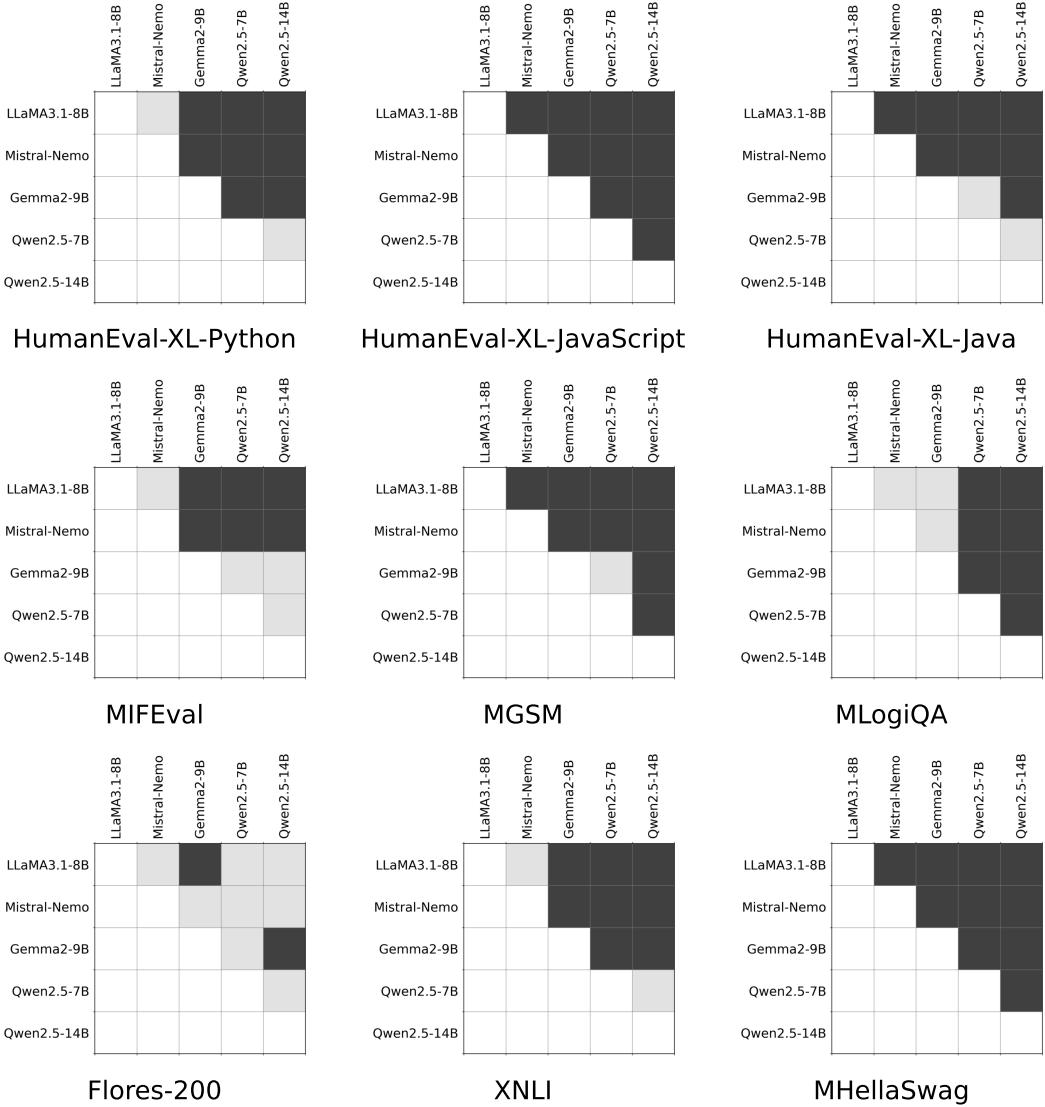


Figure 9: This figure illustrates the significant difference in pairwise performance among models with 7 to 14 billion parameters.

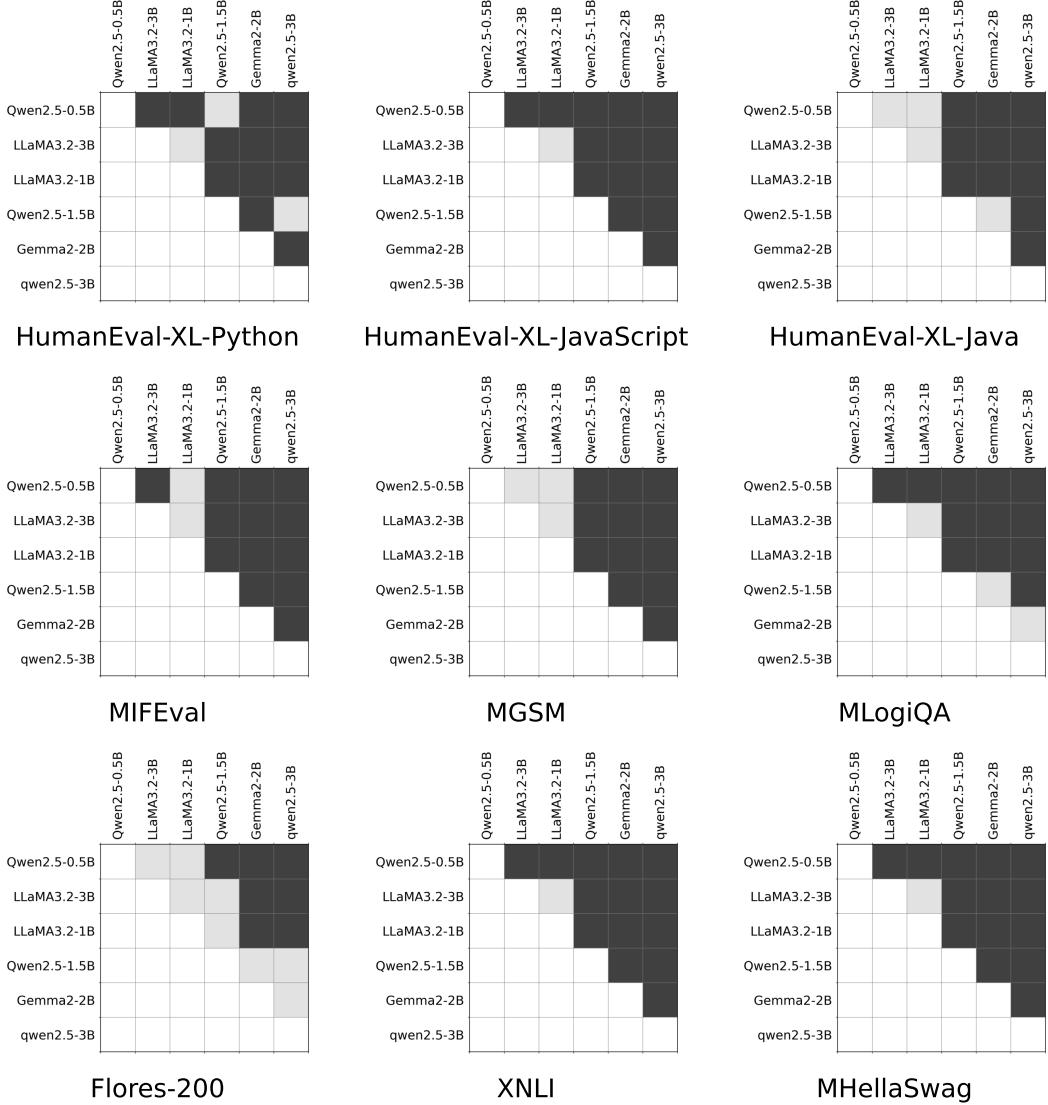


Figure 10: This figure illustrates the significant difference in pairwise performance among models with fewer than 7 billion parameters.

EN prompt for FLORES-200-en-x:

All: "Translate this sentence from English to {tgt_lang}.\n\n{src}\n"

Native prompt for FLORES-200-en-x:

zh: "将这个句子从英语翻译成中文。.\n\n{src}"

th: "แปลงประโยคนี้จากภาษาอังกฤษเป็นภาษาไทย.\n\n{src}"

ar: "قم بترجمة جملة إلى اللغة العربية.\n\n{src}"

es: "Traduce esta oración del inglés al español.\n\n{src}"

ja: "この文を英語から日本語に翻訳してください。.\n\n{src}"

ko: "이 문장을 영어에서 한국어로 번역하세요.\n\n{src}"

fr: "Traduisez cette phrase de l'anglais en français.\n\n{src}"

pt: "Traduza esta frase do inglês para o português.\n\n{src}"

vi: "Dịch câu này từ tiếng Anh sang tiếng Việt.\n\n{src}"

EN prompt for FLORES-x-en:

All: "Translate this sentence from {src_lang} to English.\n\n{src}\n"

Figure 11: This figure presents the prompt for the Flores-200 dataset.

EN prompt for MHELLASWAG:

All: "Input: {premise}\nOptions: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nPick the correct ending for the sentence from A, B, C, and D, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, C or D."

Native prompt for MHELLASWAG:

zh: "输入: {premise}\n选项: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\n从 A, B, C 或者 D 中选出正确的句子结尾，并按照以下 JSON 格式返回: \n{'answer': '[choice]'}\n其中 [choice] 必须是 A, B, C 或者 D 其中之一。"

en: "Input: {premise}\nOptions: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nPick the correct ending for the sentence from A, B, C, and D, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, C or D."

vi: "Nhập: {premise}\nLựa chọn: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nChọn kết thúc đúng cho câu từ A, B, C và D, và trả về theo định dạng JSON sau:\n{'answer': '[choice]'}\nTrong đó [choice] phải là một trong các A, B, C hoặc D."

th: "ຂ້ອມລັນເຊື່ອ: {premise}\nຕັ້ງເລືອກ: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nເລືອກຕອນຈົບທີ່ຖືກຕອງສາຫວຼນປະໄຍດວຈາ A, B, C ແລະ D ແລ້ວສັກືນໃນຮູບແບບ JSON ດັ່ງຕອບໄປນີ້:\n{'answer': '[choice]'}\nໂດຍ [choice] ຈະຕອງເປັນໜຶ່ງໃນ A, B, C ສໍາລັບ D."

ar: "أيامنلا رتختا حبل ادخال: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nنوكين أب جي ثي سنتب اعد او D، و C، و B، و A نم فلم جل حجي حصل [choice] [choice] ادح او A، و C، و B، و D."

es: "Entrada: {premise}\nOpciones: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nElija el final correcto para la oración de A, B, C y D, y devuélvalo en el siguiente formato JSON:\n{'answer': '[choice]'}\ndonde [choice] debe ser uno de A, B, C o D."

ja: "入力: {premise}\n選択肢: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nA、B、C、Dから文の正しい結末を選び、次のJSON形式で返してください: \n{'answer': '[choice]'}\nここで、[choice]はA、B、C、またはDのいずれかでなければなりません。"

ko: "입력: {premise}\n옵션: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nA, B, C, D 중에서 문장의 올바른 엔딩을 선택하고, 다음 JSON 형식으로 반환하십시오:\n{'answer': '[choice]'}\n여기서 [choice]는 A, B, C 또는 D 중 하나여야 합니다."

fr: "Entrée : {premise}\nOptions : \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nChoisissez la fin correcte de la phrase parmi A, B, C et D, et renvoyez-la dans le format JSON suivant :\n{'answer': '[choice]'}\n[choice] doit être l'un de A, B, C ou D."

pt: "Entrada: {premise}\nOpções: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nEscolha o final correto para a frase de A, B, C e D, e retorne-o no seguinte formato JSON:\n{'answer': '[choice]'}\n[choice] deve ser uma das opções A, B, C ou D."

Figure 12: This figure presents the prompt for the MHellaSwag dataset.

EN prompt for XNLI:

All: "Take the following as truth: {premise}\nThen the following statement: "{hypothesis}" is\nOptions: \nA. true\nB. inconclusive\nC. false\nSelect the correct option from A, B, and C, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, and C."

Native prompt for XNLI:

zh: "假设以下内容为真: {premise}\n考虑以下陈述: “{hypothesis}”\n该陈述是: \n选项: \nA. 真实的\nB. 无法确定\nC. 虚假的\n从 A, B 或者 C 中选择正确的选项, 并按以下JSON格式返回: \n{'answer': '[choice]'}\n其中 [choice] 必须是 A, B 或者 C 其中之一。"

en: "Take the following as truth: {premise}\nThen the following statement: "{hypothesis}" is\nOptions: \nA. true\nB. inconclusive\nC. false\nSelect the correct option from A, B, and C, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, and C."

th: "ໃຫ້ຕີວ່າເປັນຄວາມຈົງ: {premise}\nແລ້ວຂ້ອຄວາມຕ່ອໄປນີ້: "{hypothesis}" ເປັນທັວເສືອກ: \nA. ຈົງ\nB. ໄມແນນອນ\nC. ເຖິງທີ່ມເລືອກຕົວເສືອກທີ່ຄຸກຕອງຈາກ A, B, ແລະ C ແລະ ສົງເສີນໃນຮູບແບບ JSON ດັ່ງຕອໄປນີ້:\n{'answer': '[choice]'}\nໂດຍທີ່ [choice] ຕອງເປັນໜຶ່ງໃນ A, B, ແລະ C."

ar: "فَقْرِئُوكَ يَدِي ام رَبْتَعَ: {premise}\nهَذِهِ {hypothesis} مُصْلَحٌ أَمْ مُكَلَّلٌ؟\nيُثْبِتُ حَصْنَكَ بِالْمُجَلَّلِ\nأَمْ يُنْكَلِّلُ حَصْنَكَ؟\nإِنَّ {hypothesis} مُصْلَحٌ إِذَا وَجَدَتْ حَصْنَكَ مُجَلَّلًا، وَإِنَّهُ مُكَلَّلٌ إِذَا وَجَدَتْ حَصْنَكَ مُصْلَحًا.\n[choice] تَوْجِيدُهُ مُجَلَّلًا وَ[choice] نَكَلُونَهُ مُصْلَحًا."

es: "Tome lo siguiente como verdad: {premise}\nEntonces la siguiente afirmación: "{hypothesis}" es\nOpciones:\nA. verdadera\nB. inconclusa\nC. falsa\nSeleccione la opción correcta de A, B y C, y devuélvala en el siguiente formato JSON:\n{'answer': '[choice]'}\ndonde [choice] debe ser una de A, B y C."

ja: "次の内容を真実とみなしてください: {premise}\n次の文: "{hypothesis}" は\n選択肢: \nA. 真\nB. 不確定\nC. 偽\nA、B、Cの中から正しい選択肢を選び、次のJSON形式で返してください: \n{'answer': '[choice]'}\nここで、[choice]はA、B、Cのいずれかでなければなりません。"

ko: "다음 내용을 진실로 간주하십시오: {premise}\n그렇다면 다음 진술: "{hypothesis}"는\n옵션: \nA. 사실\nB. 결론을 내릴 수 없음\nC. 거짓\nA, B, C 중에서 올바른 옵션을 선택하고 다음 JSON 형식으로 반환하십시오:\n{'answer': '[choice]'}\n여기서 [choice]는 A, B 및 C 중 하나여야 합니다."

fr: "Prenez ce qui suit comme vérité : {premise}\nAlors, l'affirmation suivante : "{hypothesis}" est\nOptions : \nA. vraie\nB. inconclusive\nC. fausse\nSélectionnez l'option correcte parmi A, B et C, puis renvoyez-la dans le format JSON suivant :\n{'answer': '[choice]'}\n[choice] doit être l'un de A, B et C."

pt: "Considere o seguinte como verdade: {premise}\nEntão, a seguinte afirmação: "{hypothesis}" é\nOpções: \nA. verdadeira\nB. inconclusiva\nC. falsa\nSelecione a opção correta de A, B e C e retorne-a no seguinte formato JSON:\n{'answer': '[choice]'}\n[choice] deve ser uma das opções A, B ou C."

vi: "Xem điều sau đây là đúng: {premise}\nVậy tuyên bố sau đây: "{hypothesis}" là\nCác lựa chọn: \nA. đúng\nB. không kết luận\nC. sai\nChọn lựa chọn đúng từ A, B và C, và trả lại nó theo định dạng JSON sau:\n{'answer': '[choice]'}\ntrong đó [choice] phải là một trong A, B và C."

Figure 13: This figure presents the prompt for the XNLI dataset.

Native prompt for MGSM:

en: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "The answer is ". Do not add anything other than the integer answer after "The answer is ".\n\n{question}"

es: "Resuelve este problema matemático. Proporciona los pasos de razonamiento antes de dar la respuesta final en la última línea por sí misma en el formato de "La respuesta es ". No añadas nada más que la respuesta entera después de "La respuesta es ".\n\n{question}"

fr: "Résolvez ce problème de mathématiques. Donnez les étapes de raisonnement avant de fournir la réponse finale sur la dernière ligne elle-même dans le format de "La réponse est ". N'ajoutez rien d'autre que la réponse entière après "La réponse est ".\n\n{question}"

ja: "の数学の問題を解いてください。最終的な答えを出す前に、解答の推論過程を記述してください。
そして最後の行には"答えは"の形式で答えを記述し、その後には整数の答え以外何も追加しないでください。 \n\n{question}"

th: "แก้ปัญหาคณิตศาสตร์นี้ ให้ให้ขั้นตอนการใช้เหตุผลก่อนที่จะให้คำตอบสุดท้ายในบรรทัดสุดท้ายโดยอยู่ในรูปแบบ "คำตอบคือ" ในคราวเพิ่มของในนอกจากคำตอบที่เป็นจำนวนเต็มหลัง "คำตอบคือ"\n\n{question}"

zh: "解决这个数学问题。在最后一行给出答案前，请提供推理步骤。最后一行应该以"答案是"的形式独立给出答案。在"答案是"后不要添加除整数答案之外的任何内容。 \n\n{question}"

ar: "نقـت مـتـي نـأـبـجـي لـجـلـا تـاـوـطـخـ مـيـدـقـتـى جـرـي، رـيـخـلـا رـطـسـلـا يـفـ قـبـاجـلـا اـطـعـلـبـقـ. فـيـضـ اـيـرـلـا فـلـأـسـمـلـا ذـهـلـجـبـ مـقـ
لـا دـدـعـلـا عـوـسـ " وـهـ بـاـوـجـلـا " دـعـبـ ئـيـشـ يـأـفـيـضـتـ الـ. " وـهـ بـاـوـجـلـا " لـكـشـبـ رـيـخـلـا رـطـسـلـا يـفـ قـبـاجـلـا مـيـ
صـ حـيـحـصـ"\n\n{question}"

ko: "이 수학 문제를 해결하십시오. 마지막 줄에 답을 제시하기 전에 추론 단계를 제공하십시오. 마지막 줄은 "답변은" 형식으로 독립적으로 답을 제시해야 합니다. "답변은" 뒤에는 정수답 이외의 어떤 것도 추가하지 마십시오.\n\n{question}"

pt: "Resolva este problema matemático. Antes de dar a resposta na última linha, por favor, forneça os passos de raciocínio. A última linha deve apresentar a resposta de forma independente, começando com "A resposta é ". Após "A resposta é " não adicione nada além da resposta em número inteiro.\n\n{question}"

vi: "Giải quyết vấn đề toán học này. Trước khi đưa ra đáp án ở dòng cuối cùng, hãy cung cấp các bước lập luận. Dòng cuối cùng nên đưa ra đáp án dưới dạng "Câu trả lời là" một cách độc lập. Không thêm bất cứ nội dung nào ngoài đáp án là số nguyên sau "Câu trả lời là".\n\n{question}"

Figure 14: This figure presents the Native prompt for the MGSM dataset.

EN prompt for MGSM:

en: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "The answer is ". Do not add anything other than the integer answer after "The answer is ".\n\n{question}"

es: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "La respuesta es ". Do not add anything other than the integer answer after "La respuesta es ".\n\n{question}"

fr: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "La réponse est ". Do not add anything other than the integer answer after "La réponse est ".\n\n{question}"

ja: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "答えは ". Do not add anything other than the integer answer after "答えは ".\n\n{question}"

th: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "คําตอบคํ ". Do not add anything other than the integer answer after "คําตอบคํ ".\n\n{question}"

zh: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "答案是 ". Do not add anything other than the integer answer after "答案是 ".\n\n{question}"

ar: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "و ه ب ا و ج ل ". Do not add anything other than the integer answer after "و ه ب ا و ج ل ".\n\n{question}"

ko: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "답변은 ". Do not add anything other than the integer answer after "답변은 ".\n\n{question}"

pt: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "A resposta é ". Do not add anything other than the integer answer after "A resposta é ".\n\n{question}"

vi: "Solve this math problem. Give the reasoning steps before giving the final answer on the last line by itself in the format of "Câu trả lời là ". Do not add anything other than the integer answer after "Câu trả lời là ".\n\n{question}"

Figure 15: This figure presents the EN prompt for the MGSM dataset.

EN prompt for MLOGIQA:

All: "Passage: {context}\nQuestion: {question}\nChoices:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nPlease choose the most suitable one among A, B, C and D as the answer to this question, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, C and D."

Native prompt for MLOGIQA:

zh: "段落: {context}\n问题: {question}\n选择:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n请在 A、B、C 和 D 中选择最合适的一个作为此问题的答案，并以以下 JSON 格式返回:\n{'answer': '[choice]'}\n其中 [choice] 必须是 A、B、C 和 D 中的一项。"

en: "Passage: {context}\nQuestion: {question}\nChoices:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nPlease choose the most suitable one among A, B, C and D as the answer to this question, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, C and D."

vi: "Đoạn văn: {context}\nCâu hỏi: {question}\nLựa chọn:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nVui lòng chọn câu trả lời phù hợp nhất trong số A, B, C và D cho câu hỏi này, và trả lại nó trong định dạng JSON sau:\n{'answer': '[choice]'}\ntrong đó [choice] phải là một trong A, B, C và D."

th: "ข้อความ: {context}\nคำถาม: {question}\nตัวเลือก:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nโปรดเลือกขอที่เหมาะสมที่สุดจาก A, B, C และ D เป็นคำตอบของคำถาม และส่งคืนในรูปแบบ JSON ดังต่อไปนี้:\n{'answer': '[choice]'}\nโดยที่ [choice] จะต้องเป็นหนึ่งใน A, B, C และ D."

ar: "أهلاً وسهلاً بـ {context}\nالسؤال: {question}\nالإجابة: \nA. {option_1}\nB. {option_2}\nC. {option_3}\nD. {option_4}\nنوكري نأبجي ثيحيح كل انتلـ {question} قيسـنـتـبـ ادعـأـوـ Dـ وـ Cـ وـ Bـ وـ Aـ نـمـ فـلـمـجـلـلـ ظـحـيـ حـصـلـ [choice]ـ وـ أـ وـ Cـ وـ Aـ وـ Dـ."

es: "Pasaje: {context}\nPregunta: {question}\nOpciones:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nPor favor, elija la más adecuada entre A, B, C y D como respuesta a esta pregunta, y devuélvala en el siguiente formato JSON:\n{'answer': '[choice]'}\ndonde [choice] debe ser uno de A, B, C y D."

ja: "本文: {context}\n質問: {question}\n選択肢:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nこの質問の答えとして A、B、C、D の中から最も適したものを選択し、次の JSON 形式で返してください:\n{'answer': '[choice]'}\nここで [choice] は A、B、C、または D のいずれかでなければなりません。"

ko: "구문: {context}\n질문: {question}\n선택:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n이 질문의 답으로 A, B, C 및 D 중 가장 적합한 것을 선택하고, 다음 JSON 형식으로 반환하십시오:\n{'answer': '[choice]'}\n여기서 [choice]는 A, B, C 및 D 중 하나여야 합니다."

fr: "Passage : {context}\nQuestion : {question}\nChoix :\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nVeuillez choisir le plus approprié parmi A, B, C et D comme réponse à cette question, et le renvoyer dans le format JSON suivant :\n{'answer': '[choice]'}\non où [choice] doit être l'un de A, B, C ou D."

pt: "Passagem: {context}\nPergunta: {question}\nOpções:\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\nPor favor, escolha a mais adequada entre A, B, C e D como resposta a esta pergunta, e retorne-a no seguinte formato JSON:\n{'answer': '[choice]'}\ndonde [choice] deve ser uma das opções A, B, C ou D."

Figure 16: This figure presents the prompt for the MLogiQA dataset.

EN prompt for MMMLU:

All: "The following is a multiple-choice question. Please choose the most suitable one among A, B, C and D as the answer to this question, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, C and D.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

Native prompt for MMMLU:

zh: "以下是一个多项选择题。请在 A、B、C 和 D 中选择最合适的一个作为此问题的答案，并以以下 JSON 格式返回：\n{'answer': '[choice]'}\n其中 [choice] 必须是 A、B、C 和 D 中的一项。{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

en: "The following is a multiple-choice question. Please choose the most suitable one among A, B, C and D as the answer to this question, and return it in the following JSON format:\n{'answer': '[choice]'}\nwhere [choice] must be one of A, B, C and D.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

vi: "Đuôi đây là một câu hỏi trắc nghiệm. Vui lòng chọn câu trả lời phù hợp nhất trong số A, B, C và D cho câu hỏi này, và trả lại nó trong định dạng JSON sau:\n{'answer': '[choice]'}\ntrong đó [choice] phải là một trong A, B, C và D.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

th: "ຕ້ອໄປນີ້ສອຄໍາຄາມແບບເລືອກຕອບໜາຍຕ້ວເລືອ ໂປຣດເລືອກຂ້ອທີ່ເໝາະສົມທີ່ສຸດຈາກ A, B, C ແລະ D ເປັນຄໍາຕອບຂອງຄໍາຄາມນ ແລະສົງເສົນໃນຮູບແບບ JSON ຕ້ອໄປນີ້:\n{'answer': '[choice]'}\nໂດຍທີ່ [choice] ຈະຕອງເປັນຫຸ້ນໃນ A, B, C ແລະ D. \n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

ar: "لاب متداع او، لاؤسلا اذه ىلع قباج إك D و C و B و A نعيب نم بسن آلا راي تختا ىجري. تارا اي خل ال ددعتم لاؤس وه يلاتل ا يلاتل ا JSON نم ادح او A و B و C و D \n{'answer': '[choice]'}\n[choice] نوكي نأ بج يثي حنونت \n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

es: "Lo siguiente es una pregunta de opción múltiple. Por favor, elija la más adecuada entre A, B, C y D como respuesta a esta pregunta, y devuélvala en el siguiente formato JSON:\n{'answer': '[choice]'}\n donde [choice] debe ser uno de A, B, C y D.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

ja: "以下は選択式の質問です。この質問の答えとして A、B、C、D の中から最も適したものを選択し、次の JSON 形式で返してください: \n{'answer': '[choice]'}\nここで [choice] は A、B、C、D のいずれかでなければなりません。{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

ko: "다음은 객관식 질문입니다. 이 질문의 답으로 A, B, C 및 D 중 가장 적합한 것을 선택하고 다음 JSON 형식으로 반환하십시오: \n{'answer': '[choice]'}\n여기서 [choice]는 A, B, C 및 D 중 하나여야 합니다.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

fr: "Ce qui suit est une question à choix multiple. Veuillez choisir la plus appropriée parmi A, B, C et D comme réponse à cette question, et la renvoyer dans le format JSON suivant:\n{'answer': '[choice]'}\n où [choice] doit être l'un de A, B, C ou D.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

pt: "O seguinte é uma questão de múltipla escolha. Por favor, escolha a mais adequada entre A, B, C e D como resposta a esta pergunta, e retorne-a no seguinte formato JSON:\n{'answer': '[choice]'}\n onde [choice] deve ser uma das opções A, B, C ou D.\n\n{question}\nA. {option_a}\nB. {option_b}\nC. {option_c}\nD. {option_d}\n"

Figure 17: This figure presents the prompt for the MMMLU dataset.