

# 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. These insights offer valuable guidance for future research <sup>1</sup>.

## 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-

<sup>1</sup>We will publish all the code and resources after the paper is received.

R (Ruder et al., 2021), XGLUE (Liang et al., 2020), 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 datasets. 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) dataset for natural language inference, XCOPA (Ponti et al., 2020), MHEL-LASWAG (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. The BUFFET benchmark also provides a fixed set of few-shot demonstrations for evaluation.

**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

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)	120 × 10 (3)	Translation	Acc
MHELLASWAG (Dac Lai et al., 2023)		120 × 10 (3)	Translation	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).

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 mathematical reasoning, the OpenAI multilingual version of MMLU (MMMLU)<sup>2</sup> for knowledge comprehension, and HUMANEVAL-XL (Chen et al., 2021) for code generation.

All the benchmarks mentioned above focus either exclusively on fundamental NLP capabilities or on advanced application abilities. Additionally, there is inconsistent multilingual coverage across various datasets within a single multi-task benchmark. The proposed benchmark P-MMEVAL integrates three fundamental NLP datasets and five capability-specialized datasets, providing consistent language coverage across all selected datasets.

### 3 Datasets Selection Pipeline

Through the accumulation of a long time, the evaluation tasks for language models encompass a wide variety, with each category amassing substantial multilingual datasets. These datasets are primarily categorized into two main types: generation and understanding. Each task is further divided into various subcategories, most of which consist of multiple datasets. Therefore, selecting effective ones is crucial, as it can reduce redundant testing and improve evaluation efficiency. To achieve this, we utilize paired-sample T-test (Field, 2005) to optimize the selection process by filtering out datasets that can effectively distinguish the performances of LLMs among different model series and sizes. We

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

suggest that if these benchmarks do not maintain significant differences even when the size gap is large enough, their evaluation results can be considered ineffective. Therefore, those benchmarks can not present reliable and meaningful performance identification and comparison.

Our selection pipeline can be described as follows: Given the evaluation results of model  $A$  and model  $B$  on a multilingual dataset  $D$ , denoted as  $A_i$  and  $B_i$  respectively, where  $i$  represents the language index. Following this, we first collect two score arrays  $[A_1, A_2, \dots, A_m]$  and  $[B_1, B_2, \dots, B_m]$  which represents the evaluation results of model  $A$  and model  $B$  on  $m$  different languages, respectively. Then, we use these two arrays to derive the significance value  $p$  after running a paired-T significance test. If  $p$  is less than a pre-defined significance level (e.g., 0.01), it can be concluded that there is a significant difference in the overall scores between model  $A$  and model  $B$ . By determining whether multiple pairs of models have significantly different scores on this dataset, the effectiveness of the dataset in distinguishing the performance among various models can be identified.

### 4 P-MMEval

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

## 4.1 Design Principles

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

**Diversity in languages** To ensure that our benchmark can also help testify the cross-lingual transferability of LLMs, we unify 10 different languages spanning 8 language families, including English (*en*), Chinese (*zh*), Arabic (*ar*), Spanish (*es*), Japanese (*ja*), Korean (*ko*), Thai (*th*), French (*fr*), Portuguese (*pt*), and Vietnamese (*vi*).

## 4.2 Fundamental NLP Dataset Curation

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

### 4.2.1 Tasks

**Natural Language Understanding (NLU)** Here, we have five different sub-tasks: i) The natural language inference (NLI) dataset, XNLI (Conneau et al., 2018), which involves classifying whether a hypothesis is entailed, contradicted, or unrelated to the premise. ii) Three commonsense reasoning datasets encompass XCOPA (Ponti et al., 2020) focusing on causal reasoning, MHELLASWAG examining social scenarios and linguistic fluency, and XWINOGRAD (Tikhonov and Ryabinin, 2021) addressing anaphora resolution issues. iii) The paraphrase identification dataset PAWS-X (Yang et al., 2019) requires the model to determine whether two given sentences convey the same meaning. iv) The word sense disambiguation dataset XL-WIC (Raganato et al., 2020) focuses on understanding the meanings of words in various contexts. v) Three span-prediction datasets, i.e., XQUAD (Artetxe et al., 2020), MLQA (Lewis et al., 2020), and TYDIQA-GOLDP (Joshi et al., 2017b), where the answer to a question is provided within a piece of context.

**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.

### 4.2.2 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).<sup>3</sup> 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.

### 4.2.3 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 and extension.

## 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.<sup>4</sup> To maintain consistency across all languages, we extend the support of some benchmark datasets on the missing languages by collecting human-annotated translation results. We first deliver the translated examples

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

<sup>4</sup>For 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
XCOPIA	✗	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.

generated by powerful LLM, and require a professional translation team to conduct a thorough review of the machine translation results, correct translation errors if necessary, localize vocabulary expressions, and eliminate cases that cannot be directly mapped across languages, thus ensuring translation quality and cultural adaptability (See Table 6). In detail, the 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). We collect 80 examples in *ja*, *ko*, and *th* in extension.
- **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. We extend its multilingual support with *ar*, *ko*, *pt*, and *vi* examples.
- **Logical reasoning** We keep the original *en* and *zh* examples from origin LOGIQA (Liu et al., 2020) dataset. Besides, we extend its multilingual version by translating *en* examples into *ar*, *es*, *ja*, *ko*, *th*, *fr*, *pt*, and *vi*.
- **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. We translate those examples into *th* and *fr*.

- **Instruction following** We employ the English IFEVAL (Liu et al., 2020) dataset, which consists examples following pre-defined 25 types of “verifiable instruction”. We also extend its multilingual version MIFeVal with the support in *zh*, *ar*, *es*, *ja*, *ko*, *th*, *fr*, *pt*, and *vi*, where 96 examples for each language.

#### 4.4 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 conducted in Table 3.

Model	Understanding		Code generation	Mathematical reasoning	Logic reasoning	Knowledge	Instruction following	Generation	AVG_S	AVG_U
	XNLI	MHELLASWAG	HUMANEVAL-XL	MGSM	MLOGIQA	MMMLU	MIFEVAL	FLORES-200		
<i>Open-source models (&lt;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	<b>36.85</b>	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	<b>56.45</b>	<b>54.09</b>
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	<b>36.48</b>	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	<b>69.20</b>	<b>68.80</b>
<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	<b>71.95</b>	<b>72.36</b>
GEMMA2-27B	68.00	64.12	76.67	85.28	50.50	49.42	81.35	<b>42.23</b>	68.64	66.06
<i>Open-source models (&gt;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	<b>73.13</b>	<b>73.69</b>
MISTRAL-LARGE	69.58	69.04	77.17	90.48	53.50	51.85	83.23	<b>43.40</b>	71.25	69.31
<i>Close-source models</i>										
GPT-4o	69.17	81.04	77.05	91.60	56.75	55.77	85.21	<b>46.32</b>	73.28	<b>75.11</b>
CLAUDE-3.5-SONNET	71.50	77.72	82.92	92.84	62.25	56.17	80.73	16.20	<b>74.98</b>	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. HUMANEVAL-XL score presents the average score of three programming languages.

## 5.1 Multilingual Models

We evaluate the performance of several representative instruction-tuned models – (i) closed-source models GPT-4o<sup>5</sup> (OpenAI, 2023) and CLAUDE-3.5-SONNET<sup>6</sup>, (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.

For MGSM, we employ Chain of Thought (CoT) (Wei et al., 2022) reasoning, which guides the

<sup>5</sup>gpt-4o-2024-05-13

<sup>6</sup>claude-3-5-sonnet-20240620

model to think step-by-step before providing a final answer. For XNLI, MHELLASWAG, MLOGIQA, HUMANEVAL-XL, MIFEVAL, and FLORES-200, direct answering is utilized, which requests the model to produce answers directly. The inference methods for these datasets align with the most commonly used settings. Notably, for MMMLU, we choose the prompt template following OpenAI simple-evals repository.<sup>7</sup> Specifically, CoT reasoning exhibits a significantly higher answer extraction failure rate compared to direct answering on small-sized LLMs (i.e., the number of parameters is less than 7B), leading to poor performance. Thus, we employ a direct answering prompt for small-sized LLMs. The detailed evaluation prompts are illustrated in Appendix G.

For the few-shot demonstrations, we primarily sample demonstrations from the validation set of the original dataset. For the missing multilingual portions, we utilize GPT-4o to translate these demonstrations from English into the missing languages. Please note that the demonstrations serve only as an answer format.

## 5.3 Main Results

Table 3 presents an overview of the evaluation results. Unless otherwise noted, the standard EN prompt is applied to all datasets except FLORES-

<sup>7</sup><https://github.com/openai/simple-evals>

200, HUMANEVAL-XL, and MIFEVAL, where the Native prompt is required. More information about the prompting strategies including EN, Native, and En-Few-Shot is shown in Appendix A. The evaluation result on HUMANEVAL-XL is the average score across three programming languages including Python, JavaScript, and Java. See Appendix C for programming language evaluation details.

First, the multilingual capabilities of models become stronger as the model sizes increase (Kaplan et al., 2020). One exception is that when the size of LLAMA3.2 increases from 1B to 3B, there is a slight decline in performance. The main reason for this is that LLAMA3.2-1B and LLAMA3.2-3B exhibit poor instruction-following capabilities, leading to a higher failure rate in answer extraction and, consequently, fluctuations in the final score. As the model size increases, the improvements in various multilingual tasks show significant differences. Evaluation results on the understanding and capability-specialized tasks show significant improvement in understanding context, processing semantic information, reasoning, and special abilities, with increasing model sizes. For example, for the QWEN2.5 series, the scores on the MGSM dataset for the 0.5B and 72B models are 13.12 and 91.00, respectively. In contrast, the models’ performance on generation tasks is relatively weaker and shows slight improvement. Evaluations on the FLORES-200 datasets indicate that, despite the increase in model size, the generation capability does not improve proportionally. This may reflect the complexity of generating text that maintains logical coherence and contextual relevance, where increasing model sizes does not significantly enhance output quality.

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

## 6 Analyses

### 6.1 Analysis on Dataset Utility

The primary objective of this section is to assess the utility of each dataset within P-MMEVAL in distinguishing model performances. We divide open-sourced models into categories by two aspects: model series and model sizes. Specifically, we collect 5 categories of models from 5 model series:

- QWEN2.5: 0.5B, 1.5B, 3B, 7B, 14B, 32B, 72B
- LLAMA3.1: 8B, 70B
- LLAMA3.2: 1B, 3B
- GEMMA2: 2B, 9B, 27B
- MISTRAL: NEMO, Large

And, we divide them into three categories based on their sizes:

- Less than 7B (<7B): QWEN2.5-0.5B, QWEN2.5-1.5B, QWEN2.5-3B, LLAMA3.2-1B, LLAMA3.2-3B, GEMMA2-2B
- Between 7B and 14B (7B-14B): QWEN2.5-7B, LLAMA3.1-8B, GEMMA2-9B, MISTRAL-NEMO, QWEN2.5-14B
- Larger than 70B (>70B): LLAMA3.1-70B, QWEN2.5-72B, MISTRAL-LARGE

Table 4 shows the utility of each dataset in distinguishing the performances of paired models within the same category. The detailed method for calculating the utility of each dataset is presented in Appendix E. 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:

- LLAMA3.2-1B and LLAMA3.2-3B show no significant performance differences across almost all datasets, indicating similar multilingual capabilities. The performance differentiation of small-size models below 7B is slightly worse.

Dataset	MISTRAL	LLAMA3.2	LLAMA3.1	QWEN2.5	GEMMA2	>70B	7B-14B	<7B
FLORES-200	2/2	2/2	1/2	4/7	3/3	3/3	2/5	3/6
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 and 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.

- Compared to JavaScript and Java, most models show poor performance differentiation in Python. According to the Appendix C, 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.
- All selected datasets can distinguish between models in the majority of categories, which verifies the effectiveness of all datasets included in P-MMEVAL.

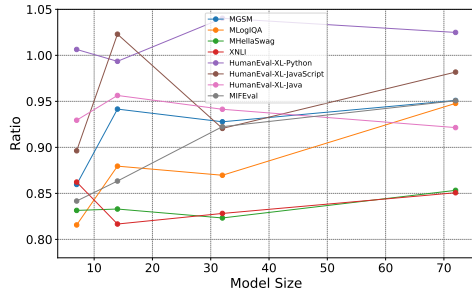


Figure 1: Illustration on the ratio of non-English performance to English performance with increasing model sizes of QWEN2.5.

## 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 use various sizes of the QWEN2.5 model (7B, 14B, 32B, and 72B) to evaluate their performance on six datasets with parallel samples in different languages. 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.

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 evaluates LLMs across 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.



621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671

## Limitations

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

1) Language Coverage: While P-MMEval currently covers 10 languages from 8 language families, there is a need to include more languages to better represent global linguistic diversity. Future work will focus on expanding the language coverage to ensure a more comprehensive evaluation of multilingual LLMs.

2) Task Diversity: P-MMEval includes 7 representative tasks, but the rapidly evolving field of LLMs demands a broader range of tasks. Future work will focus on expanding the benchmark to cover more diverse and challenging tasks, providing a more thorough assessment of multilingual LLMs.

## Ethics Statement

All procedures performed in studies involving human participants were in accordance with the ethical standards of the institutional and/or national research committee and with the 1964 Helsinki Declaration and its later amendments or comparable ethical standards. This article does not contain any studies with animals performed by any of the authors. Informed consent was obtained from all individual participants included in the study.

## References

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

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 4623–4637. Association for Computational Linguistics.

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

*Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pages 1771–1800. Association for Computational Linguistics. 672  
673  
674  
675

Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024. Mt-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 7421–7454. Association for Computational Linguistics. 676  
677  
678  
679  
680  
681  
682  
683  
684  
685

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. *Qwen technical report*. *arXiv preprint arXiv:abs/2309.16609*. 686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosiute, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemí Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. Constitutional AI: harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*. 700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*. 718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732





966	Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozinska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshv, 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> .	
967		
968		
969		
970		
971		
972		
973		
974		
975		
976		
977		
978		
979		
980		
981		
982		
983		
984	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.	
985		
986		
987		
988		
989		
990		
991		
992		
993		
994	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavathula, 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.	
995		
996		
997		
998		
999		
1000		
1001		
1002		
1003		
1004	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 multilingual chain-of-thought reasoners. In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
1005		
1006		
1007		
1008		
1009		
1010		
1011		
1012	Alexey Tikhonov and Max Ryabinin. 2021. It’s all in the heads: Using attention heads as a baseline for cross-lingual transfer in commonsense reasoning. In <i>Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021</i> , volume ACL/IJCNLP 2021 of <i>Findings of ACL</i> , pages 3534–3546. Association for Computational Linguistics.	
1013		
1014		
1015		
1016		
1017		
1018		
1019		
1020	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, Wenyin Fu, Brian Fuller,	
1021		
1022		
1023		
1024		
1025		
	Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, 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, Binh 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> .	1026
		1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
		1036
		1037
		1038
		1039
		1040
		1041
		1042
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	1043
		1044
		1045
		1046
		1047
		1048
		1049
		1050
	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. <i>Qwen2 technical report</i> . <i>arXiv preprint arXiv:abs/2407.10671</i> .	1051
		1052
		1053
		1054
		1055
		1056
		1057
		1058
		1059
		1060
		1061
		1062
		1063
		1064
		1065
		1066
		1067
	Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. PAWS-X: A cross-lingual adversarial dataset for paraphrase identification. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3685–3690. Association for Computational Linguistics.	1068
		1069
		1070
		1071
		1072
		1073
		1074
		1075
		1076
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 4791–4800. Association for Computational Linguistics.	1077
		1078
		1079
		1080
		1081
		1082
		1083
		1084

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

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

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 5: Comparison on P-MMEVAL using three different prompt settings.

## A The Impact of Different Prompts on Model Performance

We explore three different prompting strategies: EN, Native, and En-Few-Shot. Table 5 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, remaining within 2%, indicating no substantial performance 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 describe 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.

## B 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 6 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.

## C Evaluation Results on Three Programming Languages of HumanEval-XL

Table 7 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.

## D 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 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.

## E 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 (Fريتag 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  $\mathbf{m}_{\pi_j}$  in  $\Omega_i$ , it does not have a significant performance gap against any

Dataset	zh	ar	es	ja	ko	th	fr	pt	vi
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 6: The table presents the percentage of modifications made by professional translators to the machine translation results.

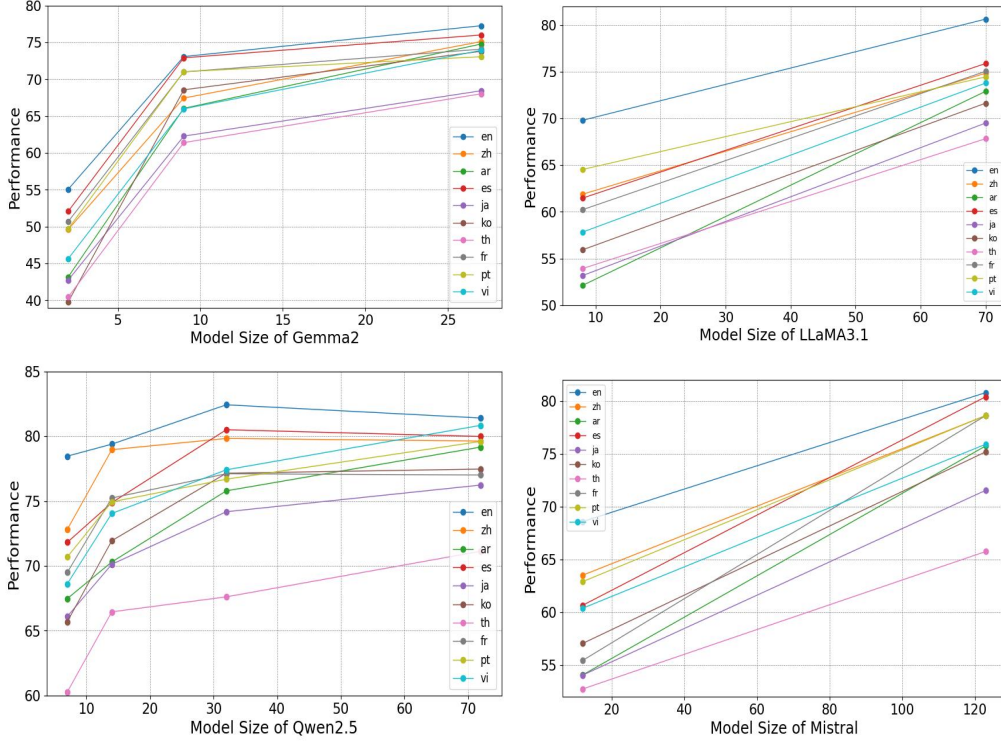


Figure 2: This figure illustrates the trend of the performance of the model in each language with increasing model sizes.

1169 model in  $\Omega_i$  except itself:

$$1170 f_1 = \begin{cases} \text{true if } \mathcal{T}(\mathbf{m}_{\pi_j}, \mathbf{m}_{\pi_p}) > \theta \text{ holds for any} \\ \quad p \in [1, 2, \dots, k], j \neq p; \\ \text{false otherwise;} \end{cases} \quad (1)$$

1171 2) for each model in  $\Omega_i$ , it has significant perfor-  
1172 mance gaps against all the model not in  $\Omega_i$ :

$$1173 f_2 = \begin{cases} \text{true if } \mathcal{T}(\mathbf{m}_{\pi_j}, \mathbf{m}_p) < \theta \text{ holds for all} \\ \quad p \notin [\pi_1, \pi_2, \dots, \pi_k]; \\ \text{false otherwise;} \end{cases} \quad (2)$$

1174 where  $\mathcal{T}(\cdot, \cdot)$  returns the  $p$ -value of the perfor-  
1175 mances between two given models, and  $\theta$  repre-  
1176 sents the threshold for denoting significance level.  
1177 The group  $\Omega_i$  is fixed if  $f_1$  and  $f_2$  both hold true.  
1178 Such a recurrent process continues till each model  
1179 is gathered into one specific group.<sup>8</sup>

1180 After gathering all models into several groups,  
1181 we use the ratio of the number of such groups to  
1182 the number of models to describe the utility of the  
1183 specific dataset. A higher ratio means that we have  
1184 more gathered groups, indicating that the bench-  
1185 mark is of high utility in distinguishing the perfor-  
1186 mances of models. On the contrary, a lower ratio

<sup>8</sup>See Algorithm 1 in Appendix E for more details.

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

1187 means that most of the models can be gathered into  
1188 the same group, denoting that the benchmark may  
1189 hardly tell which model performs better than any  
1190 other model.

1191 The algorithm for quantifying the utility of each  
1192 benchmark dataset is presented in Algorithm 1.

## 1193 **F Significance Detection on Each Dataset**

1194 The section illustrates the significant difference  
1195 between models’ pairwise performance for all cate-  
1196 gories of models.

## 1197 **G The Prompt Utilized for Each Dataset**

1198 The section presents the inference prompt utilized  
1199 for each dataset.

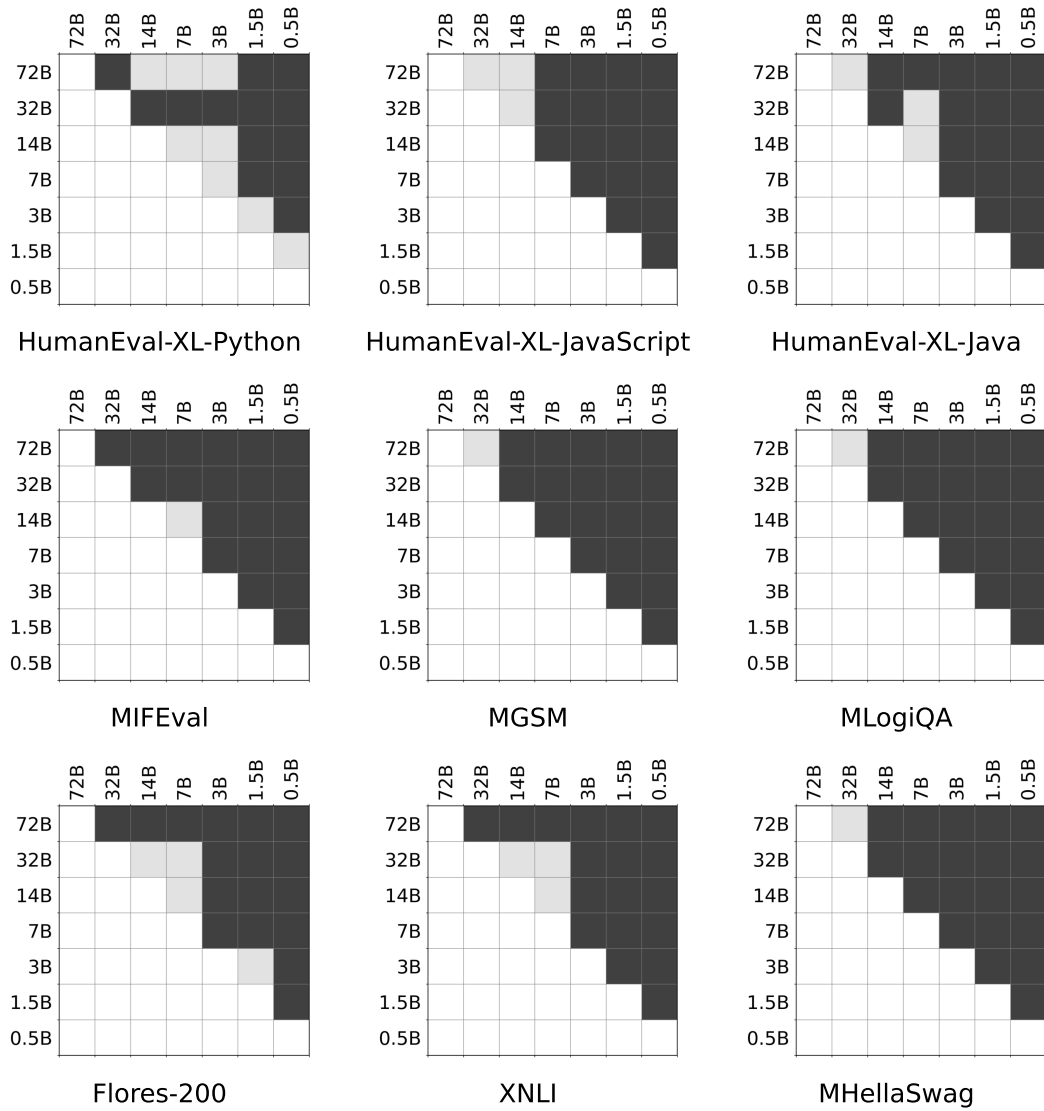


Figure 3: This figure illustrates the significant difference in pairwise performance among QWEN2.5 series models. Black blocks indicate that the  $p$ -values of paired t-tests between the corresponding models (vertical and horizontal) are less than 0.01, while gray blocks indicate  $p$ -values greater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



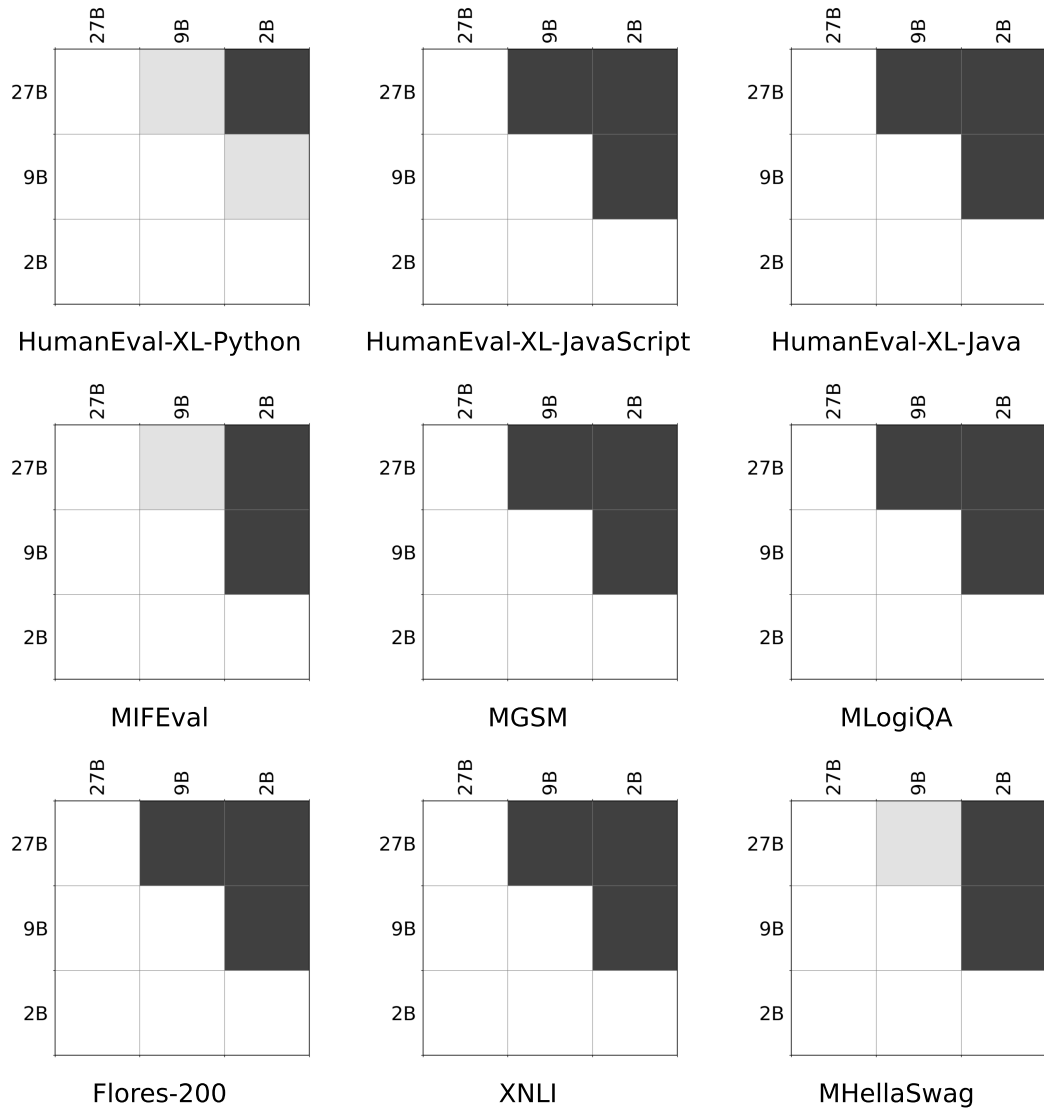


Figure 4: This figure illustrates the significant difference in pairwise performance among GEMMA2 series models.

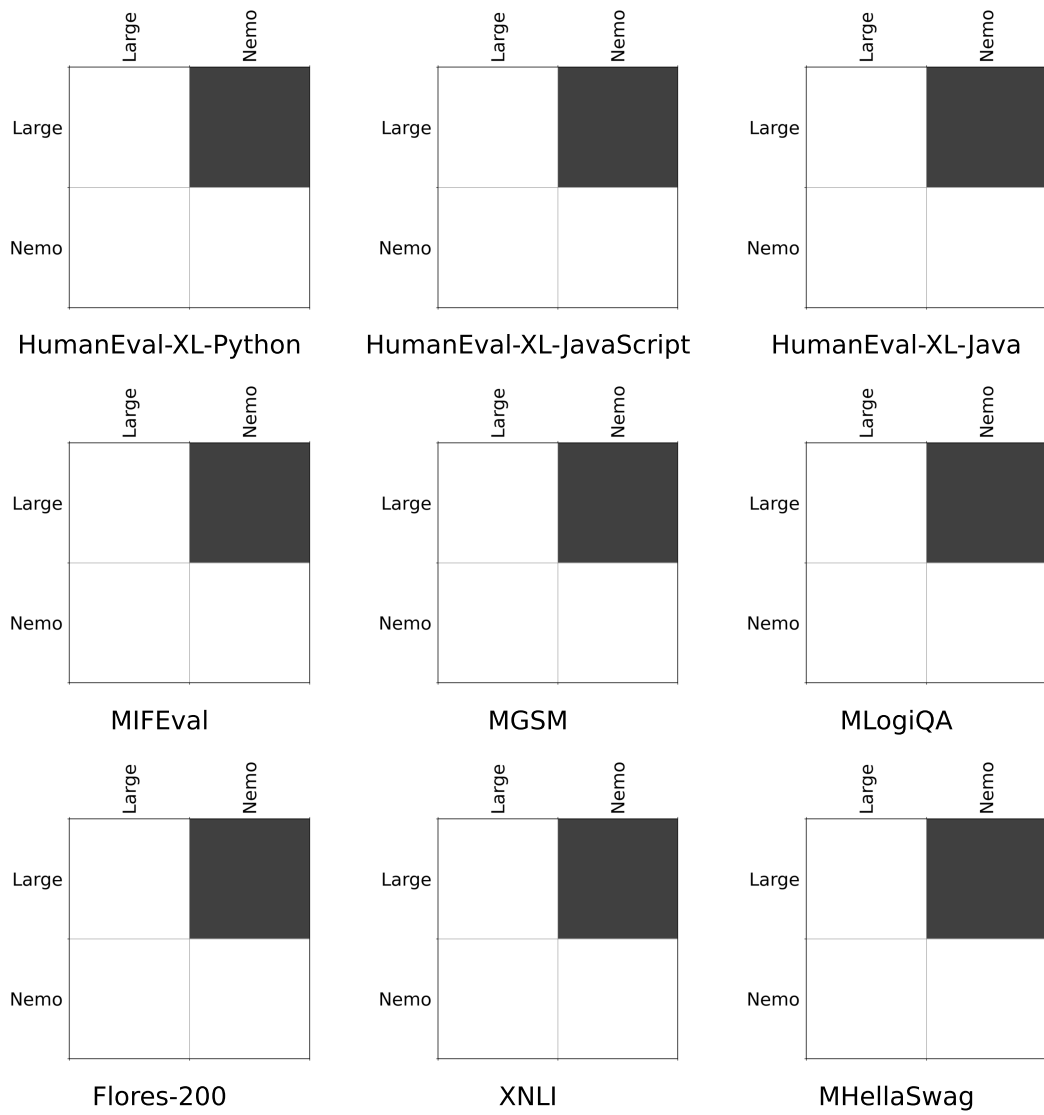


Figure 5: This figure illustrates the significant difference in pairwise performance among MISTRAL series models.

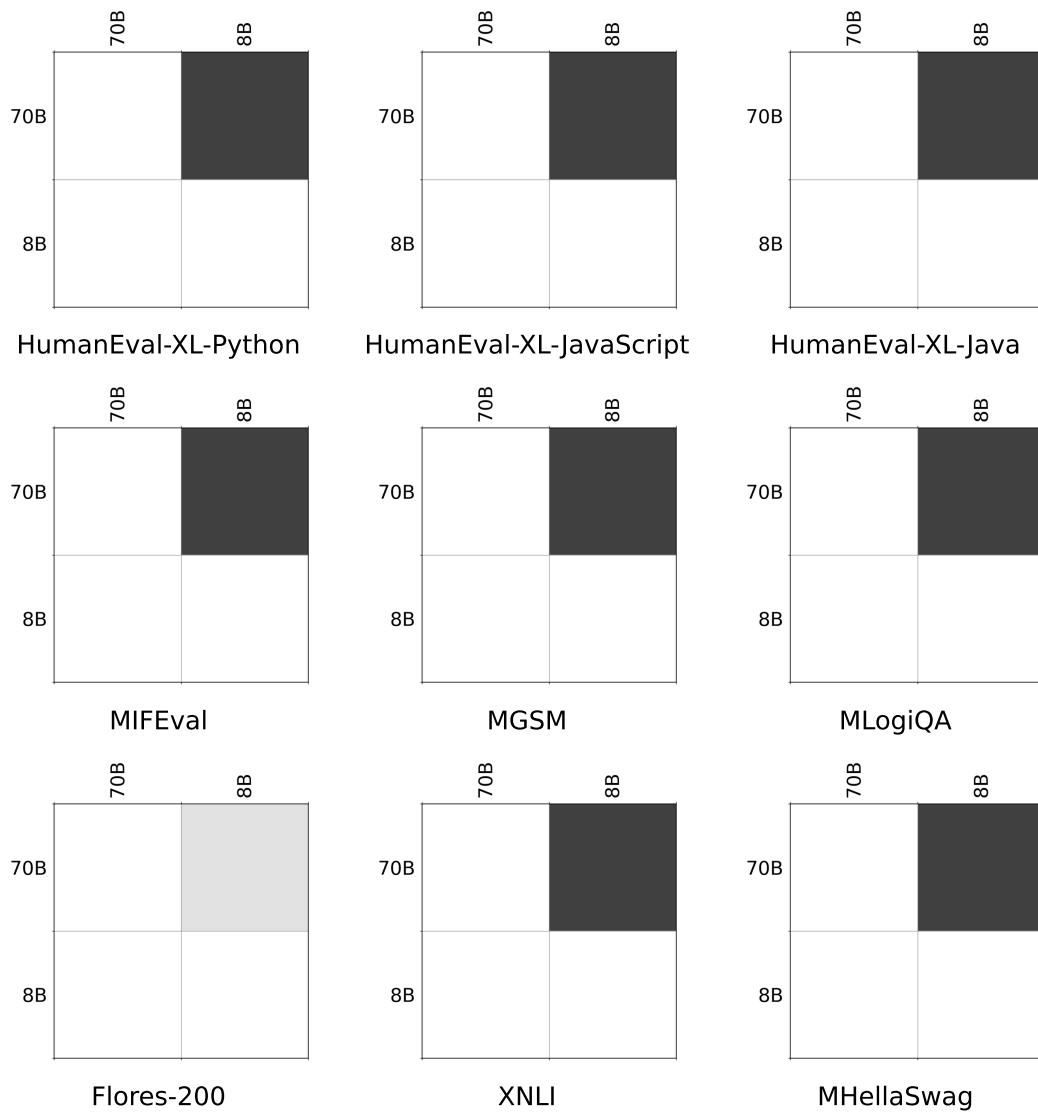


Figure 6: This figure illustrates the significant difference in pairwise performance among LLAMA3.1 series models.

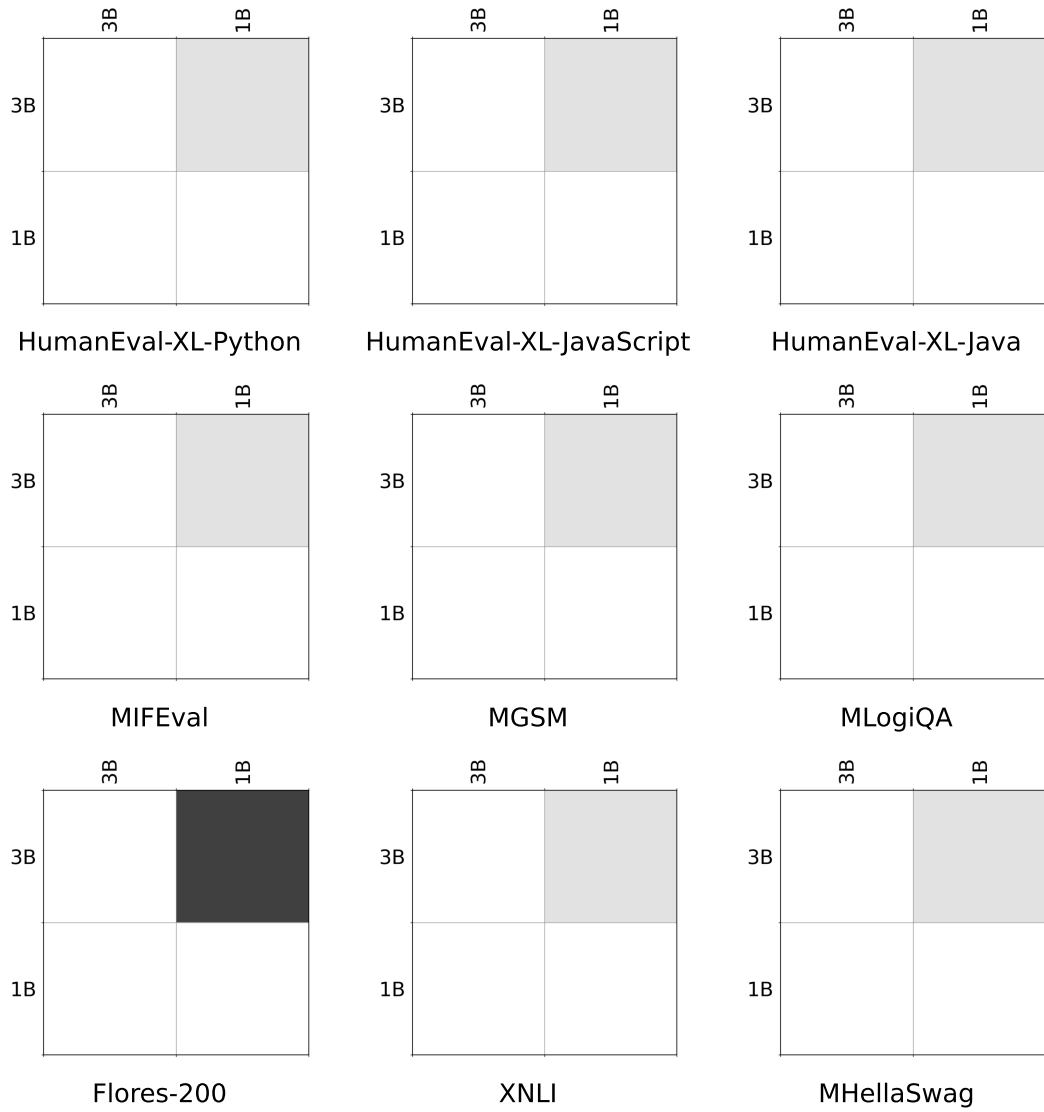


Figure 7: This figure illustrates the significant difference in pairwise performance among LLAMA3.2 series models.

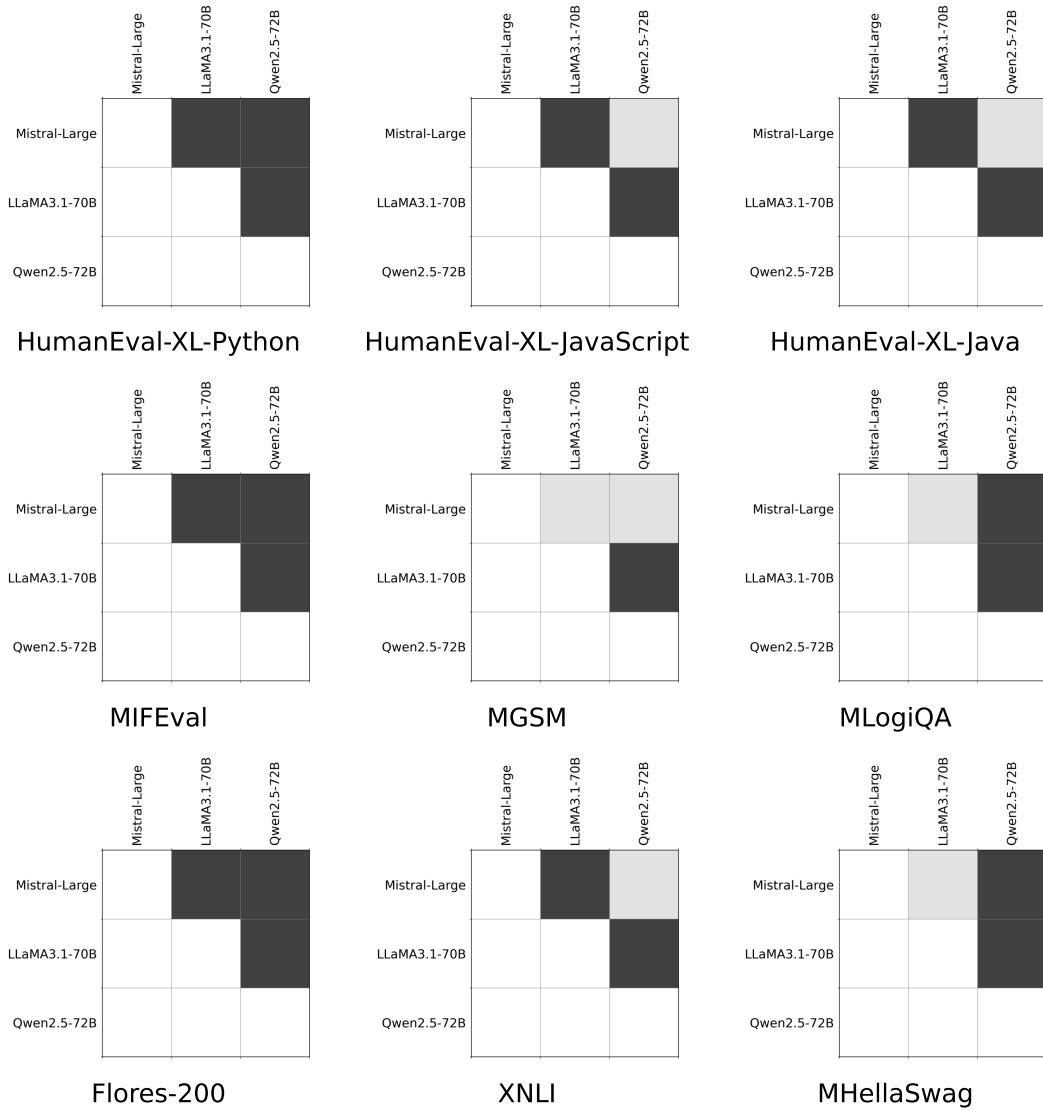


Figure 8: This figure illustrates the significant difference in pairwise performance among models with more than 70 billion parameters.

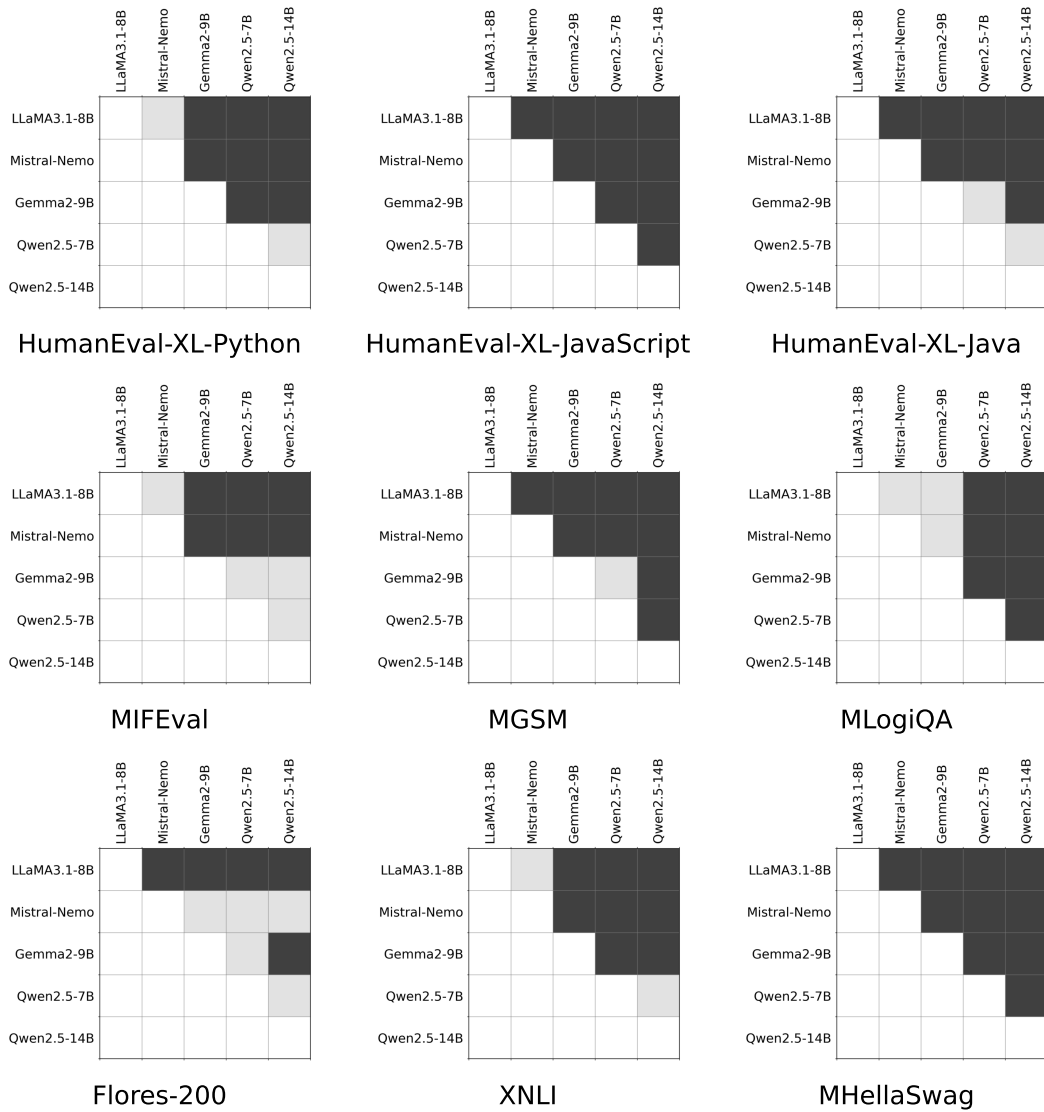


Figure 9: This figure illustrates the significant difference in pairwise performance among models with 7 to 14 billion parameters.

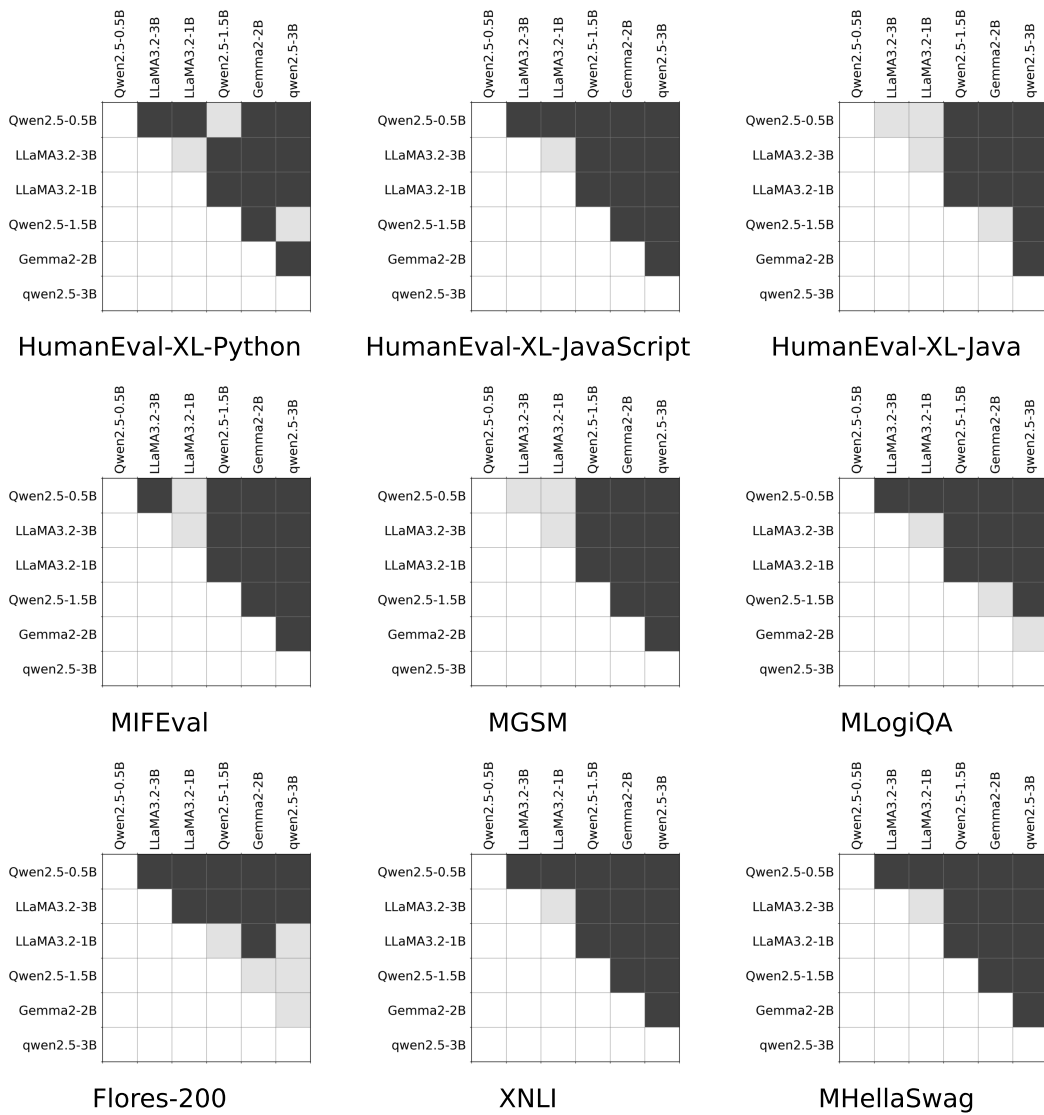


Figure 10: This figure illustrates the significant difference in pairwise performance among models with fewer than 7 billion parameters.

---

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

- 1:  $\Omega \leftarrow \emptyset$  ▷ Initialize with an empty list
- 2:  $\mathbf{z} = \mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_m$
- 3: **while**  $\mathbf{z} \neq \phi$  **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**  $\mathbf{y} \neq \phi$  **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:          $\mathbf{x} \leftarrow \mathbf{x} + \{d\}$
- 20:          $\mathbf{y} \leftarrow \mathbf{y} - \{d\}$  ▷ Moving model  $d$  into the same group
- 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

---

**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: "แปลประโยคนี้จากภาษาอังกฤษเป็นภาษาไทย.\n\n{src}"

ar: "تترجم هذا الجملة من اللغة الإنجليزية 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 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:** "ให้ถือว่าเป็นความจริง: {premise}\nแล้วข้อความต่อไปนี้: "{hypothesis}" เป็นตัวเลือก: \nA. จริง\nB. ไม่แน่นอน\nC. เท็จ\nเลือกตัวเลือกที่ถูกต้องจาก A, B, และ C และส่งคืนในรูปแบบ JSON ดังต่อไปนี้:\n{'answer': '[choice]'}\nโดยที่ [choice] ต้องเป็นหนึ่งใน A, B, และ C."

**ar:** "تتخذ التالي كحقيقة: {premise}\nثم التالي عبارة: "{hypothesis}"\nالخيارات: \nA. صحيحة\nB. غير مؤكدة\nC. خاطئة\nاختر الخيار الصحيح من A, B و C و ارجع في الصيغة JSON التالية:\n{'answer': '[choice]'}\nحيث يجب ان يكون [choice] واحد من A و B و C."

**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]'}\n donde [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 où [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 onde [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 " ไม่ควรเพิ่มอะไรนอกจากคำตอบที่เป็นจำนวนเต็มหลังจ "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:** "دقت مېتي نأ بجي .لحل ا تاواطخ مي دقت ي جري ،ري خأل رطسلا يف قبا جإا اعاطع ا لبق .ةيض اي رلا قل أسم ا هذه ل حب مق " لا ددعلا يوس " وه باوجل " دعب عيش ي أفضت ال . " وه باوجل " لكش يل ع لقتسم لكش ب ري خأل رطسلا يف قبا جإا مي باجل الل حي حص .\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 "คำตอบคือ ". Do not add anything other than the integer answer after "คำตอบคือ ".\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 "وه باوجل ". Do not add anything other than the integer answer after "وه باوجل ".\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}\nโปรดเลือกข้อที่เหมาะสมที่สุดจาก A, B, C และ D เป็นคำตอบของคำถาม และส่งคืนในรูปแบบ JSON ดังต่อไปนี้:\n{'answer': '[choice]'}\nโดยที่ [choice] จะต้องเป็นหนึ่งใน A, B, C และ D."

ar: "السياق: {context}\nالسؤال: {question}\nالخيارات:\nA. {option\_1}\nB. {option\_2}\nC. {option\_3}\nD. {option\_4}\nأخي، الرجاء اختيار الإجابة الأكثر ملاءمة من بين A، B، C و D، وارجع لي الإجابة في صيغة JSON كالتالي:\n{'answer': '[choice]'}\nحيث يجب أن تكون [choice] واحدة من A، B، C و D."

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]'}\n donde [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]'}\nou [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]'}\n onde [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。 \n\n{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:** "Dưới 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:** "ต่อไปนี้เป็นคำถามแบบเลือกตอบหลายตัวเลือก โปรดเลือกข้อที่เหมาะสมที่สุดจาก A, B, C และ D เป็นคำตอบของคำถาม และส่งคืนในรูปแบบ JSON ต่อไปนี้:\n{'answer': '[choice]'}\nโดยที่ [choice] จะต้องเป็นหนึ่งใน A, B, C และ D. \n\n{question}\nA. {option\_a}\nB. {option\_b}\nC. {option\_c}\nD. {option\_d}\n"

**ar:** "اب متداعو، لاؤسلا اذه ىلع فباكك و D و C و B و A نىب نم بسنأل راىتحا ىجرى بتارايخلا ددعتم لاؤس وه ىلاتلا لاب متداعو، لاؤسلا اذه ىلع فباكك و D و C و B و A نم أدحو [choice] نوكي نأ بجي ثي ح JSON ىلاتلا:\n{'answer': '[choice]'}\nقىسنت D.\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]'}\nndonde [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 のいずれかであればなりません。 \n\n{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]'}\nnoù [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]'}\nndonde [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.