

INCLUDE: EVALUATING MULTILINGUAL LANGUAGE UNDERSTANDING WITH REGIONAL KNOWLEDGE

Anonymous authors

Paper under double-blind review

ABSTRACT

The performance differential of large language models (LLM) between languages hinders their effective deployment in many regions, inhibiting the potential economic and societal value of generative AI tools in many communities. However, the development of functional LLMs in many languages (*i.e.*, multilingual LLMs) is bottlenecked by the lack of high-quality evaluation resources in languages other than English. Moreover, current practices in multilingual benchmark construction often translate English resources, ignoring the regional and cultural knowledge of the environments in which multilingual systems would be used. In this work, we construct an evaluation suite of 197,243 QA pairs from local exam sources to measure the capabilities of multilingual LLMs in a variety of regional contexts. Our novel resource, **INCLUDE**, is a comprehensive knowledge- and reasoning-centric benchmark across 44 written languages that evaluates multilingual LLMs for performance in the actual language environments where they would be deployed.

1 INTRODUCTION

The rapid advancement of AI technologies underscores the importance of developing LLMs that are proficient across diverse linguistic and cultural contexts, ensuring fair and equitable performance for stakeholders from various language groups. However, the lack of high-quality evaluation benchmarks in many languages discourages practitioners from training multilingual LLMs to meet this challenge. This evaluation gap limits the effective deployment of LLMs for many regions, exacerbates digital divides, and inhibits the economic and societal value of AI tools in many underserved communities.

The source of this gap is the multitude of challenges in evaluating LLMs for multilingual contexts. First, at a meta-level, the majority of benchmarks for LLMs are only in English (Hendrycks et al., 2020, *inter alia*). While non-English benchmarks exist for some tasks (Singh et al., 2024; Aakanksha et al., 2024; Pozzobon et al., 2024), they usually focus on single languages (Li et al., 2023; Koto et al., 2024), specific regions (Adelani et al., 2024; Cañete et al., 2020; Guevara-Rukoz et al., 2020; Cahyawijaya et al., 2022), or a particular domain (Wang et al., 2024a), ignoring the importance of joint evaluation to trace and unlock the benefits that multilingual capabilities could bring to low-resource languages (Pfeiffer et al., 2022; Üstün et al., 2024; Aryabumi et al., 2024).

Technical challenges also abound due to the manner in which multilingual datasets are often collected. Certain datasets are constructed using manually applied templates, resulting in low prompt and completion diversity (Muennighoff et al., 2022). Many more are composed of translations from high-resource languages (*e.g.*, English; Holtermann et al., 2024; Myung et al., 2024; Lai et al., 2023). These datasets often contain errors (Ponti et al., 2020; Plaza et al., 2024) and create *translationese artifacts* (Vanmassenhove et al., 2021; Hartung et al., 2023; Savoldi et al., 2021; Ji et al., 2023). Most importantly, they do not accurately reflect the regional and cultural contexts captured by different languages (Aakanksha et al., 2024; Awad et al., 2020; Ramezani & Xu, 2023; Singh et al., 2024). As seen in Figure 1 (a) (Regional Knowledge), a legal question posed in English, Russian, or Greek would likely reflect a user located in a different environment, where different laws may apply to respond correctly. Similarly, also seen in Figure 1 (a) (Cultural Knowledge), historical or cultural perspectives on the same topic may differ among the populaces of different regions.

To resolve this gap, we design a pipeline to collect a large multilingual language understanding benchmark (*i.e.*, **INCLUDE**) by collecting regional resources (*e.g.*, educational, professional, and practical tests) that are specific to countries and originally created by native speakers of each country’s

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

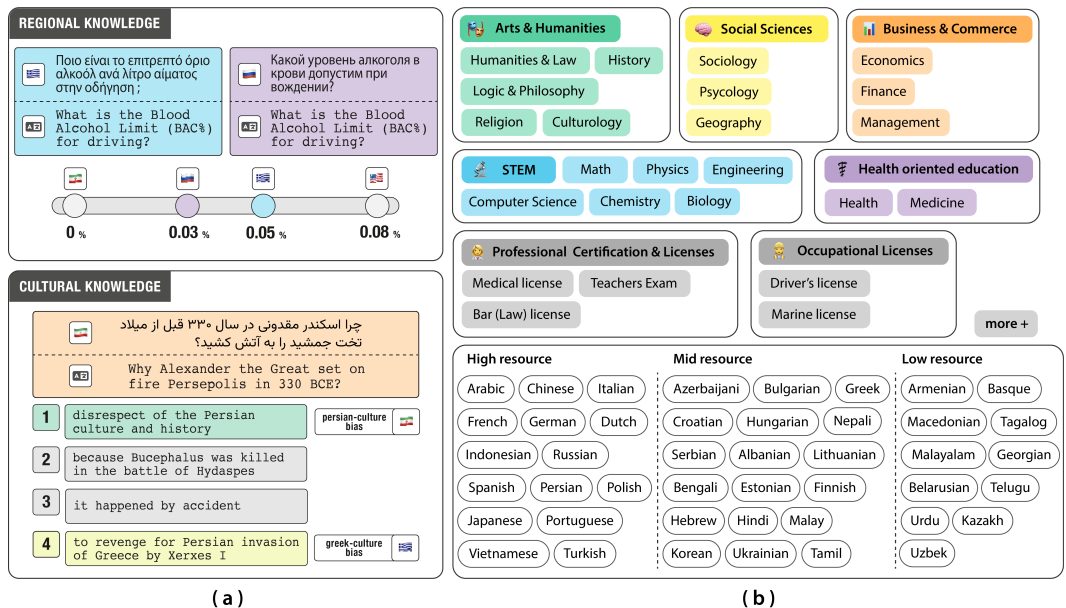


Figure 1: **Overview of INCLUDE.** (a) **Motivation:** Multilingual benchmarks must reflect the cultural and regional knowledge of the language environments in which they would be used. (b) **INCLUDE** is a multilingual benchmark compiled from academic, professional, and occupational license examinations reflecting regional and cultural knowledge in 44 languages.

official languages. This collection avoids *translationese* (Bizzoni et al., 2020) and also captures cultural nuances associated with each language, enabling rigorous evaluation of how state-of-the-art models serve diverse language users around the world.

In our experiments, we sample **INCLUDE** into two subsets for different evaluation budgets and assess an array of closed and open models on these partitions. Our results demonstrate that current models achieve high variance in performance between different languages in **INCLUDE**, and that models often struggle with questions requiring regional knowledge. Further analysis reveals that models score particularly low on languages on which they are not intentionally trained (*i.e.*, limiting regional knowledge acquisition), and that the possibility of transferring global (*i.e.*, English-aligned) perspectives improves performance for less regional topics across languages.

2 PRELIMINARIES: LANGUAGE & KNOWLEDGE

Language availability. Languages are typically characterized as a high, medium, or low resource depending on reported language availability (Joshi et al., 2020), *i.e.*, the amount of available data in a language that is available online. Interestingly, the language availability of documents used for training models (Penedo et al., 2024; Xue, 2020; Conneau et al., 2020; Computer, 2023; Üstün et al., 2024; Singh et al., 2024) differs drastically from the language distribution of non-English LLM benchmarks, with the latter being more scarce. Inspired by this discrepancy, we include 44 languages in our **INCLUDE** benchmark. In Figure 1, we characterize the availability of the included languages based on their reported availability in the mC4 corpus (Xue, 2020), and in Table 4 show further detailed metadata for each language.

Language represents regional knowledge. For LLM-based systems to be practically useful, they must enable interaction in the preferred languages of their users and be knowledgeable of the environments of those users. We define *regional knowledge* as the specific information, culture, and practices related to a local environment that is relevant for a user’s context. However, LLMs such as GPT-4 tend to exhibit a Western bias (Tao et al., 2024) due to the overrepresentation of Western text in training data (AlKhamissi et al., 2024). In **INCLUDE**, we specifically include questions encompassing the regional and cultural knowledge of a diverse set of high, medium, and low-resource languages.

3 THE INCLUDE BENCHMARK

INCLUDE is a dataset of 197,243 MCQA pairs from 1,926 examinations across 44 languages and 15 scripts. These examinations are collected from local sources in 52 countries, representing a rich array of cultural and regional knowledge. All questions in the dataset are presented in their native languages and scripts. In this section, we describe the data collection procedure for **INCLUDE**, as well as additional categorical labels we assign to each question in the dataset for later analysis.

3.1 DATA COLLECTION

To construct **INCLUDE**, we collect sources of multiple-choice exams in collaboration with native speakers and regional associations. We primarily focused on three types of exams:

Academic Exams: Exams from a variety of subjects (*e.g.*, Humanities, STEM, etc.) at different levels (*e.g.*, middle & high school, university), including country-specific national entrance exams.

Professional Certifications & Licenses: Exams issued by industry-specific regulatory bodies for specialized fields, *e.g.*, licensing exams for areas such as legal and medical practice.

Regional Licenses: Exams administered by regional authorities that assess specific qualifications, such as driving and marine licenses.

We design **INCLUDE** to assess multilingual capabilities that span beyond academic knowledge to cultural and region-specific understanding. Our data collection focuses on license and certification exams that capture regional knowledge of specific countries (in their official languages), and non-translated academic content from the humanities and social sciences to capture cultural knowledge.

From the collected sources, we extract the multiple-choice questions with their corresponding options and correct answers. More specifically, as this data came in different formats (*e.g.*, PDFs, Javascript HTML forms), we use multiple pipelines to extract QA samples from these sources and curate them in a machine-readable manner. The goal of this stage was to automate data extraction and then rely on human evaluation for verification and feature annotation.

Quality Control with Native Speakers. After automatic extraction, we provide native speakers (co-authors in this work) with parsed multiple-choice questions to ensure they were extracted correctly from source documents. In cases of extraction mistakes, annotators performed manual correction of parsed questions and answer options using the original document as a guide. In addition to performing corrections, annotators also filtered out samples that referred to images or tables, and verified that samples that rely on additional context (*e.g.*, reading comprehension) include the reference text in the question field. Finally, annotators also labeled each question with additional exam metadata, such as the language of the MCQ, its topic in both English and the original language, the academic level (if relevant), and the country of origin. In total, we parsed and verified 118,606 samples across 1,926 exam sources, amounting to 60.2% of the total data in **INCLUDE**.

Rounding out the Benchmark. To round out our benchmark, we also consolidate existing datasets with extensive domain coverage in single non-English languages: ArabicMMLU (Koto et al., 2024), ChineseMMLU (Li et al., 2023), TurkishMMLU (Yüksel et al., 2024), PersianMMLU (Ghahroodi et al., 2024) and VNHSGE (Dao et al., 2023), as well as a multilingual benchmark with limited domain coverage across multiple European languages: EXAMS (Hardalov et al., 2020). We repurpose 78,637 samples from these published benchmarks, amounting to 39.8% of the data in **INCLUDE**.

In sum, **INCLUDE** is the largest collection of multilingual exam data to-date. It is composed of 197,243 QA pairs from both novel and existing sources, and covers examinations from more than 1,926 exams in 44 different languages, 15 scripts, and 58 knowledge domains. Figure 1 and Appendix Table 4 summarize the language, domain, and knowledge diversity of the **INCLUDE** benchmark.

3.2 CATEGORIZING KNOWLEDGE

The language breadth of **INCLUDE** provides the opportunity to investigate what factors drive multilingual performance. Consequently, we annotate **INCLUDE** samples with category labels corresponding to factors such as the topic of a question and its region-specificity. Given the prohibitive cost of

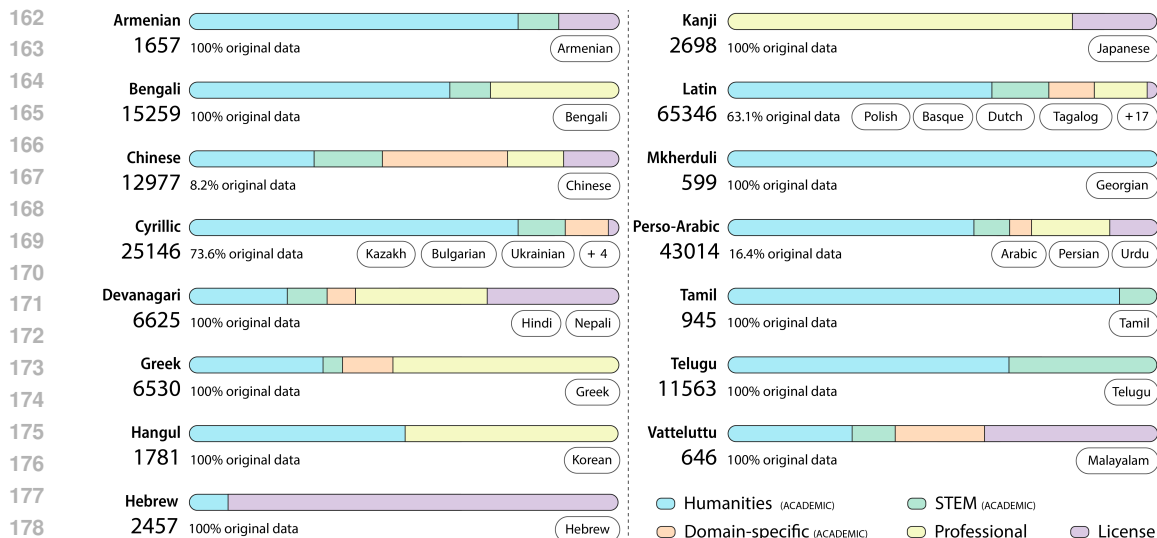


Figure 2: **Overview of the collected data grouped by script.** We depict the languages associated with each script, the total samples in each script, and the percentage of the samples that were collected from new sources that have not been published by the community yet.

performing sample-level annotation, we only perform a coarse annotation by labeling the exam sources of QA pairs, rather than individual samples. We describe our categorization schemes below.

Academic Domain. We manually categorized 1,926 unique exams, following the methodology in Hendrycks et al. (2020). Our categorization follows a two-level taxonomy: a high-level academic area (e.g., Humanities), and particular academic field within this area (e.g., History, Philosophy, Literature).¹ Each exam is categorized based on its title, which indicates its topic (e.g., Greek History) and associated level (e.g., high school, undergraduate, professional certification). Figure 7 provides a breakdown of the number of exam samples per language, organized by this taxonomy.

Regionality. To account for regional knowledge, we categorize exam questions into two major groups: **region-agnostic** and **region-specific** knowledge. Region-agnostic questions do not require knowledge of particular regions (e.g., mathematics, physics), and their answers should remain common regardless of the language in which a question is posed. In total, 34.4% of all questions collected were classified as region-agnostic. In contrast, region-specific questions require knowledge that may depend on a particular cultural or geographical context. This category is further divided into three sub-categories:

Explicitly Regional: A question is classified as *region-explicit* when it pertains to legal, regulatory, or procedural knowledge of regions. Examples include questions about local laws, certifications, clinical guidelines, or licensing requirements (see Figure 1(a)). 18.8% of all questions were *region-explicit*.

Cultural: Language often serves as an implicit marker of culture. For instance, in Figure 1(a), the answer to a question about historical figures in a Greek exam may reflect a different perspective than a similar question posed for a Persian exam. We categorize questions as *cultural* when they pertain to a region’s cultural or historical context. This category includes questions for subjects inherently tied to a region’s language, history, or social norms. 16.4% of questions were classified as *cultural*.

Implicitly Regional: Finally, the *region-implicit* category is a catch-all for other questions whose answers may depend on a certain degree of regional knowledge understanding. These questions are not explicitly regional or culture-related, but may require regional context to answer correctly. For example, business practices may be different depending on region, even if the underlying theory is common in many places. In total, 30.4% of all questions collected were classified as region-implicit.

Detailed annotation procedures for these categories are described in Appendix A.4, and general statistics about regional labels per academic area and academic field are provided in Figure 7.

¹This taxonomy is adapted from the *Outline of Academic Disciplines* found on Wikipedia.

4 EXPERIMENTAL SETUP

In this section, we describe our experimental settings for evaluating models on **INCLUDE**.

4.1 DATA SELECTION

The breadth of **INCLUDE** (197,243 QA pairs in 44 languages) makes it amenable to many evaluation use cases, including monolingual evaluation in 44 languages. However, for multilingual evaluation, this same scale is prohibitively expensive for many researchers.² Consequently, we curate two subsets of **INCLUDE** for benchmarking multilingual LLMs in different resource settings.

INCLUDE-LITE: This subset uniformly samples 22,635 QA pairs ($\sim 12\%$ of **INCLUDE**) across languages, knowledge tasks, and academic levels. The goal of this subset is to develop a multilingual benchmark with broad language and task coverage. Each language has a maximum of 550 samples, with 500 drawn from domains that correspond to regional knowledge and 50 from STEM subjects.

INCLUDE-TINY: A lightweight subset, uniformly drawn from **INCLUDE-LITE**, designed for rapid assessment of multilingual LLMs with a total of 10,770 samples ($\sim 6\%$ of **INCLUDE**). The upper limit per language is 250 samples and only includes region-specific domains.

For standardization (and alignment with prior benchmarks; Hendrycks et al., 2020), **INCLUDE-LITE** and **INCLUDE-TINY** contain only multiple-choice questions with four answer options. Questions from **INCLUDE** with fewer than four options were omitted during sampling, and questions with more than four options were pruned of options until only four remained. In the following sections, we benchmark models (§5.1) and perform analysis (§5.2-5.3) on **INCLUDE-LITE** and **INCLUDE-TINY**.³

4.2 MODELS

We assess **INCLUDE** on GPT-4o (Achiam et al., 2023) as a state-of-the-art multilingual and general-purpose model. We also investigate the role of scaling by benchmarking models that self-report parameters: we compare the larger Aya-23-35B (Aryabumi et al., 2024) 35-billion parameter model and the Qwen2.5-14B (Yang et al., 2024) 14-billion parameter model with Aya-23-8B (Aryabumi et al., 2024) and Qwen2.5-7B (with 8 and 7 billion parameters). Additionally we benchmark Mistral-7B (Jiang et al., 2023), Llama-3-7B (Dubey et al., 2024), Gemma-7B (Team et al., 2024), BLOOMZ-7B (Muennighoff et al., 2022), and XGLM-7.5B (Lin et al., 2021).

We note that some of the models we evaluate such as Mistral and Gemma, do not explicitly claim to support multiple languages, though in practice, they are heavily adopted in multilingual use cases relative to explicitly multilingual models like BLOOMZ (Muennighoff et al., 2022). Furthermore, even reportedly multilingual models (*e.g.*, Aya-23, which supports 23 languages), do not support all 44 languages included in our benchmark. Among the models, XGLM-7.5B has the widest (reported) language coverage, with 23 languages overlapping with our dataset. We evaluate all the mentioned models on the **INCLUDE-TINY** and **INCLUDE-LITE** benchmarks.

We follow the prompting strategy of Hendrycks et al. (2020) and report both 5-shot and zero-shot scores. For the zero-shot setting, we employ a Chain-of-Thought (CoT; Wei et al., 2022) approach by appending the translation of “let’s think step by step” to the prompt (Kojima et al., 2022). The maximum generation lengths for the 5-shot and zero-shot CoT configurations are set to 40 and 1024 tokens for the smaller models and to 512 and 1024 tokens for the larger models, respectively.

5 RESULTS & ANALYSIS

5.1 GENERAL PERFORMANCE

Table 1 shows the performance of all models evaluated across the 44 languages in **INCLUDE-LITE**. For larger models, *e.g.*, GPT-4o and Aya-23-35b, we provide results for both 5-shot and zero-shot CoT, while for the remaining models, we report only the 5-shot accuracy.

²The cost of evaluating **INCLUDE** using GPT-4o with 5-shot demonstrations exceeded \$1000.

³We will publicly release both **INCLUDE-LITE** and **INCLUDE-TINY**.

Model	# Langs	INCLUDE-TINY			INCLUDE-LITE		
		Total Acc.	Answer Acc.	Format Errors (%)	Total Acc.	Answer Acc.	Format Errors (%)
GPT-4o	-						
- 5-shot		77.5	82.7	6.3	77.3	83.2	7.1
- Zero-shot CoT		78.8	79.0	0.2	79.0	79.2	0.2
C4AI-Aya-23-35b	21						
- 5-shot		51.8	54.5	4.9	51.9	54.4	4.7
- Zero-shot CoT		43.3	49.3	12.3	45.7	50.0	11.1

Aya-23-8B	21	32.5	42.1	22.8	34.0	42.6	20.1
Mistral-7B (v0.3)	-	44.1	44.1	0.0	43.3	43.3	0.0
Mistral-7B-Instruct (v0.3)	-	43.8	44.0	0.3	43.6	43.8	0.4
Gemma-7B	-	55.1	55.1	0.0	54.5	54.5	0.0
Gemma-7B-Instruct	-	39.1	39.1	0.0	38.7	38.7	0.0
Qwen2.5-7B	22	53.8	54.7	1.6	54.1	55.1	1.9
Qwen2.5-7B-Instruct	22	53.2	53.4	0.4	53.8	54.0	0.5
Qwen2.5-14B	22	60.9	61.8	1.4	61.4	62.4	1.5
Llama-3-8B	-	50.0	50.0	0.0	50.3	50.3	0.0
Llama-3-8B-Instruct	-	49.9	49.9	0.0	49.8	49.8	0.0
XGLM-7.5B	23	25.3	26.4	4.0	25.0	26.1	4.1
BLOOM-7.1B	20	25.8	26.6	2.8	25.4	26.1	2.8
BLOOMZ-7.1B	20	27.8	28.8	3.4	27.8	28.8	3.6

Table 1: Results on **INCLUDE-TINY** and **INCLUDE-LITE**. **Total Accuracy** represents the raw accuracy of the model for answering **INCLUDE** questions in each respective subset. **Answer Accuracy** represents the accuracy of the model when only considering samples where an answer is extracted from the model’s output in the correct response format. **Formatting Errors (%)** describes the percentage of model responses that are not formatted correctly and so do not output any answer option. We mark these incorrect by default in **Total Accuracy** and do not include them when computing **Answer Accuracy**. **# Langs** reports the number of languages from **INCLUDE** publicly reported to be *intentionally* included in the pretraining data of each model.

Among all models, GPT-4o achieves the highest performance, reaching an accuracy of $\sim 79\%$ across all domains and examples. We observe that CoT prompting moderately enhances GPT-4o’s performance, particularly in Professional and STEM-related exams (Table 2), where the most substantial improvements were seen. In contrast, the smallest gains were observed in exams related to Licenses and the Humanities. Drawing on prior studies that compare CoT and non-CoT prompting strategies across different domains (Sprague et al., 2024), we hypothesize that this observation is due to reasoning skills required in professional examinations (e.g., medicine, law), and computation-heavy subjects in STEM. In contrast, we observe a $\sim 6\%$ performance drop for Aya-23-35b using CoT prompting on **INCLUDE-LITE**, likely because Aya-23-35b is less adapted for mathematical reasoning. GPT-4o and Aya also yield more format errors in their answers compared to other models (§5.4).

When comparing smaller models ($\leq 8\text{B}$ parameters), Gemma-7B delivers the best overall performance, with Qwen2.5-7B and Qwen2.5-7B-Instruct closely behind. While Gemma-7B excels in the Humanities and Licenses categories, the Qwen models surpass others in STEM, Applied Sciences, and Professional domains (Table 2). The lowest-performing models in our evaluation are BLOOM, BLOOMZ, and XGLM, which exhibit close to random performance. This poor performance on **INCLUDE** aligns with their results on both MMLU and translated versions of MMLU, as reported in previous studies (Ruan et al., 2024; Lai et al., 2023).

When comparing models of different sizes, we observe that the Aya-23-35B model outperforms the 8B model by $\sim 17\%$, while the Qwen2.5-14B model shows a $\sim 7\%$ improvement over its 7B counterpart on **INCLUDE-LITE**. As the pretraining data remained consistent across the different sizes within both the Aya-23 and Qwen2.5 model families, we conclude that, with similar training data, increasing model size significantly enhances multilingual capabilities.

Interestingly, we see little benefit to instruction-tuning for improving performance on **INCLUDE-LITE**. Most instruction-tuned models perform slightly worse or on par with their base counterpart, with an outlier performance drop of $\sim 15\%$ for the Gemma-7B-Instruct model. A possible explanation for

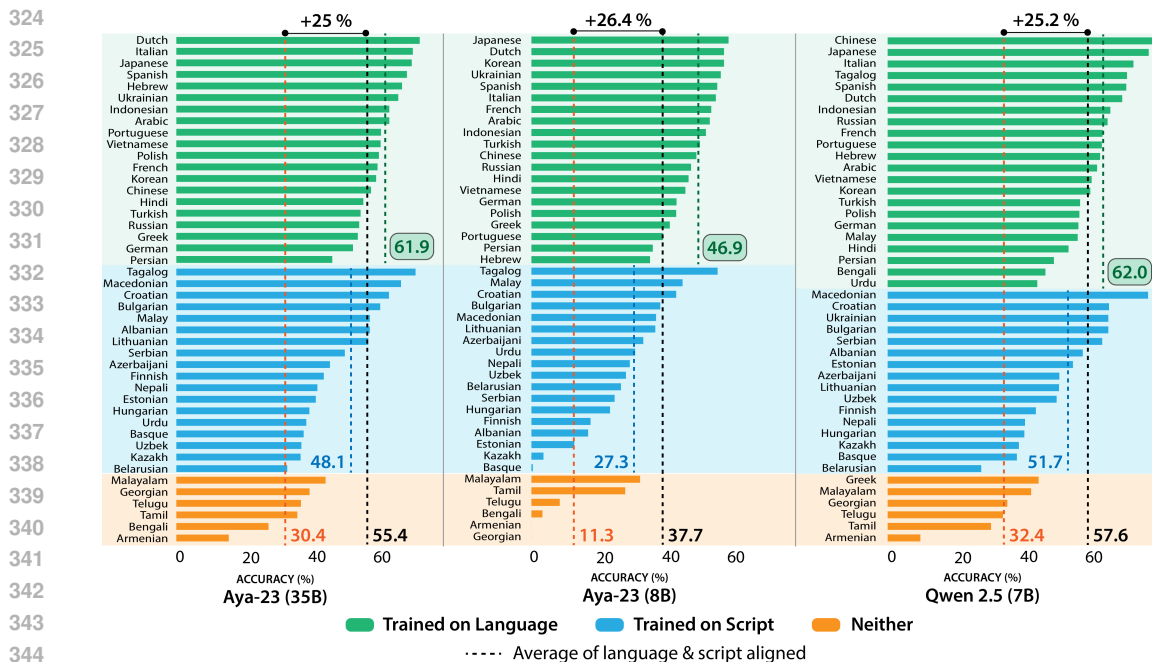


Figure 3: **Performance of models stratified by language.** Results are grouped by whether the language was explicitly included in the pretraining dataset of the model (**Trained on Language**), whether a similar language with the same script was in the pretraining corpus (**Trained on Script**), or whether there was no linguistically similar language in the pretraining corpus (**Neither**). Color dotted lines represent average performance for each category for a particular model. Black dotted lines represent average performance across all script-aligned languages.

this gap is that instruction-tuned models may have been fine-tuned predominantly on English data, potentially diminishing multilingual capabilities acquired during pretraining.

Finally, we observe performance on **INCLUDE-TINY** is roughly equivalent to **INCLUDE-LITE** for all models (within 1%), supporting its applicability for more resource-constrained evaluation settings.

5.2 LANGUAGE ANALYSIS

To better understand how LLMs perform on questions in languages seen and unseen during pretraining, we take a deeper look into three open models, *i.e.*, Aya-23-8B, Aya-23-35B, and Qwen2.5-7B, for which we have details surrounding pretraining data (and its associated language distribution). In this analysis, we specifically test three language exposure scenarios: performance on languages the model has been *intentionally*⁴ trained on, performance on languages the model was not reported to be trained on but for which the corresponding script was reported to be trained on, and performance on completely unseen languages and scripts during pretraining.

Figure 3 presents the language-stratified performance of these models on **INCLUDE-LITE**. As expected, the models demonstrate better performance on languages that were reported as part of their pretraining data (**Trained on Language**). All models also demonstrate some degree of knowledge transfer to languages they were not trained on but which share the same script as languages in their pretraining data (**Trained on Script**). In this scenario, Aya-23-35B achieves 48.1% accuracy, while Qwen2.5-7B reaches 51.7% accuracy, aligning with previous research that suggests shared scripts enable cross-lingual transfer between languages (Muller et al., 2021; Xhelili et al., 2024). Other factors may also contribute to each model’s performance on unseen languages, though, such as cross-lingual transfer across topologically-similar languages. For example, the presence of Turkish data may enhance the model’s performance on Azerbaijani (Senel et al., 2024). Pretraining data

⁴We denote *intentionally* in this context to mean that the authors specifically reported this language as being covered in the pretraining corpus of the model.

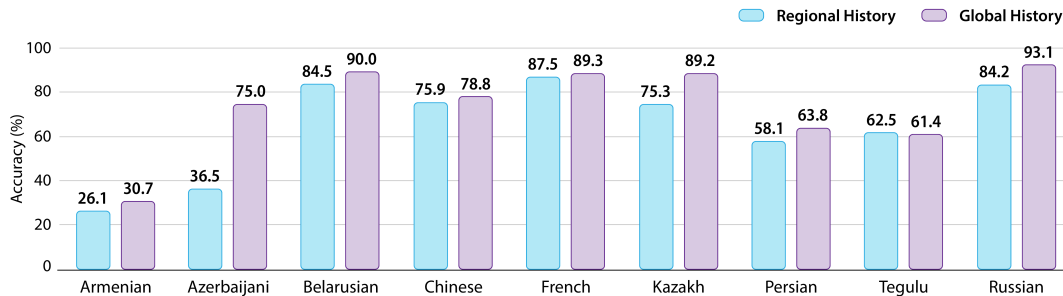


Figure 4: **GPT-4o performance** on regional history exams (*cultural*) and global history exams from that region (*region-implicit*) based on a total of 11,148 questions from **INCLUDE**. In each language (except Telugu), models perform better on the global history exam than the regional history exam.

contamination, where languages that were not intended to be in the pretraining data may still be unintentionally included (Blevins & Zettlemoyer, 2022), may also contribute to these transfer results.

Lastly, we observe that all three models perform poorly on languages whose scripts were not represented in the pretraining corpus (**Neither** in Figure 3). Data contamination is also less of a confounding factor in these cases as language identification is more robust for unique scripts (Kargaran et al., 2023). In these cases, the model often performs worse than random, likely due to not being able to produce responses in the correct format (further details in 5.4).

5.3 REGIONAL & ACADEMIC DOMAIN KNOWLEDGE PERFORMANCE

Using the category labels outlined in Section 3.2, we conduct a stratified analysis of five-shot GPT-4o’s ability to answer different types of regional questions in **INCLUDE-LITE**. We first note that overall performance differs strongly between languages. In some languages, the model consistently performs well across all academic domains and regional knowledge types, while in others, it struggles across the board (see Appendix Tables 5, 6, and 7).

However, lower performance in certain languages is often linked to questions requiring regional knowledge (*e.g.*, historical knowledge, professional certifications, medical licenses), suggesting that the model’s knowledge about regions varies significantly and that its performance across languages reflects this differential. Among regional categories, in Appendix Figure 6, we see the model performs worst on *cultural* questions, followed by *region-explicit* questions. Professional certification exams in different regions are a particular challenge for GPT-4o (average 68.6%). In Persian, the model’s accuracy on certification exams is notably low (43.2%), whereas it performs better on subjects such as Geography and Sociology (over 66%). Similarly, in Greek, GPT-4o achieves an average accuracy of 71.3%, but only 54.1% on medical license questions. Further language performance comparisons across various academic disciplines are shown in Appendix Figure 5.

Despite these findings, we note that performance across regional question types (*e.g.*, region-agnostic, cultural) is not deconfounded from other features such as topical difficulty and academic level. Indeed, we observe that *region-agnostic* questions are among the most challenging as models struggle with the mathematical nature of many STEM topics (Frieder et al., 2023; Borges et al., 2024). On average, subjects such as Mathematics and Chemistry show the lowest average accuracies (Appendix Figure 6). Unfortunately, as the region label of any question depends on its subject, it is naturally confounded with the model’s ability in that subject, regardless of whether the subject is regional or not.

History is one of the few fields where we can achieve a more controlled study of regional difference as exams can be divided into two categories: those testing region-specific historical knowledge (*e.g.*, “Armenian history”, which we label as *cultural*) and general history taught in a particular region (*e.g.*, “World History”; *region-implicit*⁵). In Figure 4, we observe that for all languages that have History exams with both *cultural* and *region-implicit* labels (with the exception of Telugu), the model

⁵We note “World History” as *region-implicit* because the manner in which the subject is taught and evaluated may vary between regions, even if the subject material seems like it should be universal.

432 performs better on the general history exams, indicating a lack of *cultural* knowledge necessary to
 433 answer questions for more region-specific topics.

434 Overall, the variance in performance among different regional categories in our results suggests that
 435 model performance on **INCLUDE** may not be rooted in across-the-board language comprehension
 436 issues, but instead in grasping specialized regional knowledge for different languages.
 437

438 5.4 CHALLENGES IN MULTILINGUAL EVALUATION 439

440 Multilingual LLMs do not generate outputs the same way in all languages. Throughout our ex-
 441 periments, we observed that models did not always follow the exact format primed by the 5-shot
 442 examples (**Format Errors** in Table 1), which required generating a longer output length to rectify.
 443 To empirically measure the impact of this seemingly minute evaluation design choice, we assess the
 444 five-shot performance of GPT-4o on **INCLUDE-LITE** across various decoded output lengths, focusing
 445 specifically on its ability to generate a correct response within the first k follow-up tokens ($k = 50$,
 446 100, 200, and 512). As in our main results, we use the 5-shot prompt template from Hendrycks et al.
 447 (2020), without explicitly instructing the model how to generate a correct answers. Instead, the model
 448 must induce the format from the provided demonstrations.

449 Table 9 presents the performance of GPT-4o across the 44 languages of the benchmark, evaluated
 450 under four different generation window settings. On average, the model shows a 3.1% performance
 451 improvement when increasing the generation length window from 50 to 512 tokens. However, this
 452 effect is not uniform; some languages experience significant improvements, such as Uzbek (+17.2%),
 453 Armenian (+13.1%), and Malayalam (+12.9%). Many others remain largely unaffected. A manual
 454 review and analysis of the generated outputs in languages with the largest gains reveal that the model
 455 often generates verbose responses, explaining the context before providing the final answer (*i.e.*,
 456 ignoring the formatting in the demonstrations, but reaching the correct response). One possible
 457 explanation for these discrepancies is the model’s limited ability to leverage in-context learning
 458 effectively in certain languages, potentially due to imbalances in language resources during the
 459 alignment phase (Zhang et al., 2024b).

460 Overall, these results demonstrate that standardizing evaluation is a challenge in tasks that may lead
 461 to different output patterns (Nayab et al., 2024), which is compounded in multilingual evaluations.
 462 In particular, given the incentive to lower generation lengths at test time (to lower inference or API
 463 costs), reliable multilingual assessment requires anticipating how models will produce outputs in
 464 different languages and scripts, and how evaluation settings might inadvertently affect measures of
 465 performance. Specifically, practitioners should be reflective about penalizing models for format errors
 466 when assessing capabilities, and intentionally probe for format errors given they may not speak or
 467 read the languages being evaluated.

468 6 RELATED WORK 469

470 In recent years, the creation of benchmarks has substantially improved the evaluation of LLMs.
 471 Pioneering efforts such as GLUE and SuperGLUE (Wang et al., 2018; 2019) played an important
 472 role in advancing tasks related to language understanding. Recent benchmarks, such as MMLU
 473 (Hendrycks et al., 2020), HellaSwag (Zellers et al., 2019), ARC (Clark et al., 2018), GSM8K (Cobbe
 474 et al., 2021), and BigBench (Srivastava et al., 2022), focus on evaluating models for more complex
 475 knowledge comprehension and reasoning. In addition to being used as final evaluations, they are
 476 often used to monitor and compare LLM performance during pretraining, rather than more traditional
 477 measures such as perplexity (Penedo et al., 2024). However, these benchmarks evaluate models only
 478 using English data, limiting their use in the development of multilingual LLMs.

479 Evaluating multilingual models requires benchmarks that assess models for these same complex
 480 abilities across diverse languages. However, initial multilingual benchmarks focus on more basic
 481 linguistic abilities (Conneau et al., 2018; Ponti et al., 2020), and collections of such tasks (Liang
 482 et al., 2020; Hu et al., 2020; Ruder et al., 2021; Asai et al., 2023; Ahuja et al., 2023a;b). Furthermore,
 483 these benchmarks generally include only a few high-resource languages or are based on translations
 484 from high-resource languages, limiting the assessment of regional knowledge comprehension and
 485 reasoning capabilities. Finally, similar to English evaluations, multilingual benchmarks have trended
 toward saturation (Zhang et al., 2024a; Wang et al., 2024b). Although there have been efforts to

486 create language-specific MMLU-like datasets, coverage remains limited to a few languages (Li
487 et al., 2023; Koto et al., 2024; Ghahroodi et al., 2024). Most similar to our proposed effort, the
488 Exams dataset (Hardalov et al., 2020) encompasses questions covering 16 languages across 24 topics
489 collected from elementary and high school science curricula. The Aya dataset (Singh et al., 2024)
490 also includes a substantial release, covering 513 million data points across 101 languages, including
491 in-language evaluation sets developed by native speakers assessing general performance and safety.
492 However, the Aya dataset is not focused on collecting in-language exams. Our work develops
493 a multilingual benchmark encompassing 44 languages, integrating questions from academic and
494 professional examinations and broadening the evaluation spectrum of multilingual LLMs to include
495 region-specific knowledge.

496 Finally, a rich body of work has developed benchmarks to assess LLMs for cultural understand-
497 ing. Arora et al. (2024) evaluate various aspects of culture and language using questions from
498 community forums on 15 topics. Aakanksha et al. (2024) curate a safety dataset that encompasses
499 local nuances. Myung et al. (2024) compile questions about food, sports, holidays, education, and
500 family translated into multiple languages. Synthetic benchmarks, such as NormAd (Rao et al., 2024),
501 generate culturally-rooted stories to measure how well models grasp societal norms. Tools such as
502 CultureBank source cultural descriptions from online platforms such as TikTok (Shi et al., 2024),
503 offering alternative ways to ground cultural benchmarks in dynamic, real-world knowledge. Beyond
504 benchmarking, Chiu et al. (2024) proposed a tool that facilitates human-machine collaboration for
505 co-creation of complex datasets, challenging the multicultural understanding and adaptability of
506 LLMs. In contrast to this line of work, our study goes beyond culture as a dimension of regional
507 knowledge, and also assesses LLMs on questions that reflect region-related factual knowledge (*e.g.*,
508 professional standards, law, clinical guidelines).

509 7 CONCLUSION

510
511 We release **INCLUDE**, a comprehensive multilingual evaluation suite designed to assess performance
512 of large language models (LLMs) across a wide range of subjects and languages for a rich array of
513 cultural and regional knowledge. **INCLUDE** contains 197,243 MCQA pairs from 1,926 examinations
514 across 44 languages and 15 scripts collected from 52 countries. Overall, our results from evaluating 15
515 models on **INCLUDE** indicate there remains considerable room for model improvement in multilingual
516 regional knowledge understanding and that regional knowledge understanding varies significantly
517 across languages. **INCLUDE** offers researchers and developers a novel and valuable benchmark for
518 evaluating and improving the regional understanding abilities of future multilingual models in the
519 language environments where they would be used.

520 8 LIMITATIONS

521
522 Our work has several limitations. First, the benchmark spans 44 languages with varying levels of
523 resource availability, leading to different distributions of questions from various academic disciplines
524 across languages. This disparity makes it challenging to perform direct comparisons between
525 performance in disciplines across languages. Additionally, the difficulty of exams may vary not
526 only between languages but also within the same language if exams originate from different sources.
527 However, this limitation is also a reflection of one of the strengths of our benchmark. Questions are
528 sourced from local examinations that reflect the regional and cultural nuances of the environments
529 in which those exams are implemented, which was our motivation for a new evaluation benchmark.
530 Naturally, this precludes exact correspondence between questions across languages. Another practical
531 limitation is that our regional knowledge labels were annotated at the exam topic level, rather than at
532 the individual question level. As a result, questions are classified based on their overarching topic,
533 rather than individual content.

ETHICS STATEMENT

The primary goal of our benchmark is to reduce disparities in performance regional knowledge understanding across languages, addressing the inequities in access to technology and its benefits that often result from these gaps. We have designed the benchmark to reflect a diverse range of linguistic, cultural, and regional contexts, sourcing data from local and region-specific exam materials. Throughout the data collection process, we ensured that no private or sensitive information was included. We only collected data from exams for which there were no license issues. Our benchmark aims to capture and integrate essential cultural knowledge across many languages. We emphasize the importance of local engagement and encourage developers using this benchmark for the evaluation of monolingual models to actively consult with local stakeholders. To promote equitable access to technology and the development of multilingual large language models, we release our benchmark to the community.

REPRODUCIBILITY STATEMENT

We plan to release two subsets, **INCLUDE-LITE** and **INCLUDE-TINY**, alongside the associated documentation and code for data processing and evaluation. These resources will be made publicly available upon acceptance. To mitigate the risk of data contamination during fine-tuning, **INCLUDE** will be released in incremental batches over a period of four months.

Further details regarding experimental settings, including resource utilization, hyperparameters, and baseline configurations, can be found in Section 4 and Appendix A.5, which provide a comprehensive overview of our methodology.

REFERENCES

- Aakanksha, Arash Ahmadian, Beyza Ermis, Seraphina Goldfarb-Tarrant, Julia Kreutzer, Marzieh Fadaee, and Sara Hooker. The multilingual alignment prism: Aligning global and local preferences to reduce harm, 2024. URL <https://arxiv.org/abs/2406.18682>.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- David Ifeoluwa Adelani, Jessica Ojo, Israel Abebe Azime, Jian Yun Zhuang, Jesujoba O. Alabi, Xuanli He, Millicent Ochieng, Sara Hooker, Andiswa Bukula, En-Shiun Annie Lee, Chiamaka Chukwunke, Happy Buzaaba, Blessing Sibanda, Godson Kalipe, Jonathan Mukiibi, Salomon Kabongo, Foutse Yuehgoh, Mmasibidi Setaka, Lolwethu Ndolela, Nkiruka Odu, Rooweither Mabuya, Shamsuddeen Hassan Muhammad, Salomey Osei, Sokhar Samb, Tadesse Kebede Guge, and Pontus Stenetorp. Irokobench: A new benchmark for african languages in the age of large language models, 2024. URL <https://arxiv.org/abs/2406.03368>.
- Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, et al. Mega: Multilingual evaluation of generative ai. *arXiv preprint arXiv:2303.12528*, 2023a.
- Sanchit Ahuja, Divyanshu Aggarwal, Varun Gumma, Ishaan Watts, Ashutosh Sathe, Millicent Ochieng, Rishav Hada, Prachi Jain, Maxamed Axmed, Kalika Bali, et al. Megaverse: Benchmarking large language models across languages, modalities, models and tasks. *arXiv preprint arXiv:2311.07463*, 2023b.
- Badr AlKhamissi, Muhammad ElNokrashy, Mai AlKhamissi, and Mona Diab. Investigating cultural alignment of large language models. *arXiv preprint arXiv:2402.13231*, 2024.
- Shane Arora, Marzena Karpinska, Hung-Ting Chen, Ipsita Bhattacharjee, Mohit Iyyer, and Eunsol Choi. Calmqa: Exploring culturally specific long-form question answering across 23 languages, 2024. URL <https://arxiv.org/abs/2406.17761>.

- 594 Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat
595 Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo,
596 Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh
597 Fadaee, Ahmet Üstün, and Sara Hooker. Aya 23: Open weight releases to further multilingual
598 progress. 2024. URL <https://arxiv.org/abs/2405.15032>.
- 599
600 Akari Asai, Sneha Kudugunta, Xinyan Velocity Yu, Terra Blevins, Hila Gonen, Machel Reid, Yulia
601 Tsvetkov, Sebastian Ruder, and Hannaneh Hajishirzi. Buffet: Benchmarking large language models
602 for few-shot cross-lingual transfer. *arXiv preprint arXiv:2305.14857*, 2023.
- 603
604 Edmond Awad, Sohan Dsouza, Azim Shariff, Iyad Rahwan, and Jean-François Bonnefon. Universals
605 and variations in moral decisions made in 42 countries by 70,000 participants. *Proceedings of the
606 National Academy of Sciences*, 117(5):2332–2337, 2020.
- 607
608 Yuri Bizzoni, Tom S Juzek, Cristina España-Bonet, Koel Dutta Chowdhury, Josef van Genabith, and
609 Elke Teich. How human is machine translation? comparing human and machine translations
610 of text and speech. In Marcello Federico, Alex Waibel, Kevin Knight, Satoshi Nakamura, Her-
611 mann Ney, Jan Niehues, Sebastian Stüker, Dekai Wu, Joseph Mariani, and Francois Yvon (eds.),
612 *Proceedings of the 17th International Conference on Spoken Language Translation*, pp. 280–290,
613 Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.iwslt-1.34.
614 URL <https://aclanthology.org/2020.iwslt-1.34>.
- 615
616 Terra Blevins and Luke Zettlemoyer. Language contamination helps explains the cross-lingual
617 capabilities of English pretrained models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang
618 (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*,
619 pp. 3563–3574, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
620 Linguistics. doi: 10.18653/v1/2022.emnlp-main.233. URL [https://aclanthology.org/
2022.emnlp-main.233](https://aclanthology.org/2022.emnlp-main.233).
- 621
622 Beatriz Borges, Negar Foroutan, Deniz Bayazit, Anna Sotnikova, Syrielle Montariol, Tanya Nazaret-
623 zky, Mohammadreza Banaei, Alireza Sakhaeirad, Philippe Servant, Seyed Parsa Neshaei, Jibril
624 Frej, Angelika Romanou, Gail Weiss, Sepideh Mamooler, Zeming Chen, Simin Fan, Silin Gao,
625 Mete Ismayilzada, Debjit Paul, Alexandre Schöpfer, Andrej Janchevski, Anja Tiede, Clarence
626 Linden, Emanuele Troiani, Francesco Salvi, Freya Behrens, Giacomo Orsi, Giovanni Piccioli,
627 Hadrien Sevel, Louis Coulon, Manuela Pineros-Rodriguez, Marin Bonnassies, Pierre Hellich, Puck
628 van Gerwen, Sankalp Gambhir, Solal Pirelli, Thomas Blanchard, Timothée Callens, Toni Abi
629 Aoun, Yannick Calvino Alonso, Yuri Cho, Alberto Chiappa, Antonio Sclocchi, Étienne Bruno,
630 Florian Hofhammer, Gabriel Pescia, Geovani Rizk, Leello Dadi, Lucas Stoffl, Manoel Horta
631 Ribeiro, Matthieu Bovel, Yueyang Pan, Aleksandra Radenovic, Alexandre Alahi, Alexander
632 Mathis, Anne-Florence Bitbol, Boi Faltings, Cécile Hébert, Devis Tuia, François Maréchal, George
633 Candea, Giuseppe Carleo, Jean-Cédric Chappelier, Nicolas Flammarion, Jean-Marie Fürbringer,
634 Jean-Philippe Pellet, Karl Aberer, Lenka Zdeborová, Marcel Salathé, Martin Jaggi, Martin Rajman,
635 Mathias Payer, Matthieu Wyart, Michael Gastpar, Michele Ceriotti, Ola Svensson, Olivier Lévêque,
636 Paolo lenne, Rachid Guerraoui, Robert West, Sanidhya Kashyap, Valerio Piazza, Viesturs Simanis,
637 Viktor Kuncak, Volkan Cevher, Philippe Schwaller, Sacha Friedli, Patrick Jermann, Tanja Kaser,
638 and Antoine Bosselut. Could chatgpt get an engineering degree? evaluating higher education
639 vulnerability to ai assistants, 2024. URL <https://arxiv.org/abs/2408.11841>.
- 640
641 Samuel Cahyawijaya, Holy Lovenia, Alham Fikri Aji, Genta Indra Winata, Bryan Wilie, Rahmad
642 Mahendra, Christian Wibisono, Ade Romadhony, Karissa Vincentio, Fajri Koto, et al. Nusacrowd:
643 Open source initiative for indonesian nlp resources. *arXiv preprint arXiv:2212.09648*, 2022.
- 644
645 José Cañete, Gabriel Chaperon, Rodrigo Fuentes, Jou-Hui Ho, Hojin Kang, and Jorge Pérez. Spanish
646 pre-trained bert model and evaluation data. In *PMLADC at ICLR 2020*, 2020.
- 647
648 Yu Ying Chiu, Liwei Jiang, Maria Antoniak, Chan Young Park, Shuyue Stella Li, Mehar Bhatia,
649 Sahithya Ravi, Yulia Tsvetkov, Vered Shwartz, and Yejin Choi. Culturalteaming: Ai-assisted
650 interactive red-teaming for challenging llms’ (lack of) multicultural knowledge, 2024. URL
<https://arxiv.org/abs/2404.06664>.

- 648 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
649 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
650 *arXiv preprint arXiv:1803.05457*, 2018.
- 651 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
652 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
653 math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>, 2021.
- 654 Together Computer. Redpajama: an open dataset for training large language models, October 2023.
655 URL <https://github.com/togethercomputer/RedPajama-Data>.
- 656 Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R Bowman, Holger
657 Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence representations. *arXiv*
658 *preprint arXiv:1809.05053*, 2018.
- 659 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek,
660 Francisco Guzm'an, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Un-
661 supervised cross-lingual representation learning at scale. In *Proceedings of the 58th An-
662 nual Meeting of the Association for Computational Linguistics*, pp. 8440–8451, Online, July
663 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747. URL
664 <https://www.aclweb.org/anthology/2020.acl-main.747>.
- 665 Xuan-Quy Dao, Ngoc-Bich Le, The-Duy Vo, Xuan-Dung Phan, Bac-Bien Ngo, Van-Tien Nguyen,
666 Thi-My-Thanh Nguyen, and Hong-Phuoc Nguyen. Vnhsge: Vietnamese high school graduation
667 examination dataset for large language models. *arXiv preprint arXiv:2305.12199*, 2023.
- 668 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
669 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
670 *arXiv preprint arXiv:2407.21783*, 2024.
- 671 Simon Frieder, Luca Pinchetti, Alexis Chevalier, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas
672 Lukaszewicz, Philipp Christian Petersen, and Julius Berner. Mathematical capabilities of chatgpt,
673 2023. URL <https://arxiv.org/abs/2301.13867>.
- 674 Omid Ghahroodi, Marzia Nouri, Mohammad Vali Sanian, Alireza Sahebi, Doratossadat Dast-
675 gheib, Ehsaneddin Asgari, Mahdieh Soleymani Baghshah, and Mohammad Hossein Rohban.
676 Khayyam challenge (persianmmlu): Is your llm truly wise to the persian language? *arXiv preprint*
677 *arXiv:2404.06644*, 2024.
- 678 Adriana Guevara-Rukoz, Isin Demirsahin, Fei He, Shan-Hui Cathy Chu, Supheakmungkol Sarin,
679 Knot Pipatsrisawat, Alexander Gutkin, Alena Butryna, and Oddur Kjartansson. Crowdsourcing
680 latin american spanish for low-resource text-to-speech. In *Proceedings of the Twelfth Language*
681 *Resources and Evaluation Conference*, pp. 6504–6513, 2020.
- 682 Momchil Hardalov, Todor Mihaylov, Dimitrina Zlatkova, Yoan Dinkov, Ivan Koychev, and Preslav
683 Nakov. Exams: A multi-subject high school examinations dataset for cross-lingual and multilingual
684 question answering. *arXiv preprint arXiv:2011.03080*, 2020.
- 685 Kai Hartung, Aaricia Herygers, Shubham Vijay Kurlekar, Khabbab Zakaria, Taylan Volkan, Sören
686 Gröttrup, and Munir Georges. Measuring sentiment bias in machine translation. In *International*
687 *Conference on Text, Speech, and Dialogue*, pp. 82–93. Springer, 2023.
- 688 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
689 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
690 *arXiv:2009.03300*, 2020.
- 691 Carolin Holtermann, Paul Röttger, Timm Dill, and Anne Lauscher. Evaluating the elementary
692 multilingual capabilities of large language models with multiq. *arXiv preprint arXiv:2403.03814*,
693 2024.
- 694 Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson.
695 Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation.
696 In *International Conference on Machine Learning*, pp. 4411–4421. PMLR, 2020.

- 702 Meng Ji, Pierrette Bouillon, and Mark Seligman. *Translation Technology in Accessible Health*
703 *Communication*. Cambridge University Press, 2023.
704
- 705 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
706 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
707 *Mistral 7b*. *arXiv preprint arXiv:2310.06825*, 2023.
- 708 Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and
709 fate of linguistic diversity and inclusion in the nlp world. *arXiv preprint arXiv:2004.09095*, 2020.
710
- 711 Amir Hossein Kargaran, Ayyoob Imani, François Yvon, and Hinrich Schuetze. GlotLID: Language
712 identification for low-resource languages. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.),
713 *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 6155–6218, Sing-
714 apore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
715 findings-emnlp.410. URL [https://aclanthology.org/2023.findings-emnlp.](https://aclanthology.org/2023.findings-emnlp.410)
716 410.
- 717 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
718 language models are zero-shot reasoners. In *Neural Information Processing Systems*, 2022. URL
719 <https://api.semanticscholar.org/CorpusID:249017743>.
720
- 721 Fajri Koto, Haonan Li, Sara Shatnawi, Jad Doughman, Abdelrahman Boda Sadallah, Aisha Alraeesi,
722 Khalid Almubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, et al. Arabicmmlu: Assessing
723 massive multitask language understanding in arabic. *arXiv preprint arXiv:2402.12840*, 2024.
- 724 Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan Rossi, and Thien
725 Nguyen. Okapi: Instruction-tuned large language models in multiple languages with reinforcement
726 learning from human feedback. In Yansong Feng and Els Lefever (eds.), *Proceedings of the 2023*
727 *Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp.
728 318–327, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/
729 v1/2023.emnlp-demo.28. URL <https://aclanthology.org/2023.emnlp-demo.28>.
730
- 731 Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy
732 Baldwin. Cmmmlu: Measuring massive multitask language understanding in chinese. *arXiv preprint*
733 *arXiv:2306.09212*, 2023.
- 734 Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou,
735 Daxin Jiang, Guihong Cao, et al. Xglue: A new benchmark dataset for cross-lingual pre-training,
736 understanding and generation. *arXiv preprint arXiv:2004.01401*, 2020.
737
- 738 Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle
739 Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. Few-shot learning with multilingual language
740 models. *arXiv preprint arXiv:2112.10668*, 2021.
- 741 Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le
742 Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. Crosslingual
743 generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*, 2022.
744
- 745 Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, and Djamé Seddah. When being unseen
746 from mBERT is just the beginning: Handling new languages with multilingual language models.
747 In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven
748 Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021*
749 *Conference of the North American Chapter of the Association for Computational Linguistics:*
750 *Human Language Technologies*, pp. 448–462, Online, June 2021. Association for Computational
751 Linguistics. doi: 10.18653/v1/2021.naacl-main.38. URL [https://aclanthology.org/](https://aclanthology.org/2021.naacl-main.38)
752 2021.naacl-main.38.
- 753 Junho Myung, Nayeon Lee, Yi Zhou, Jiho Jin, Rifki Afina Putri, Dimosthenis Antypas, Hsuvas
754 Borkakoty, Eunsu Kim, Carla Perez-Almendros, Abinew Ali Ayele, et al. Blend: A benchmark for
755 llms on everyday knowledge in diverse cultures and languages. *arXiv preprint arXiv:2406.09948*,
2024.

- 756 Sania Nayab, Giulio Rossolini, Giorgio Buttazzo, Nicolamaria Manes, and Fabrizio Giacomelli. Con-
757 cise thoughts: Impact of output length on llm reasoning and cost. *arXiv preprint arXiv:2407.19825*,
758 2024.
- 759
760 Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro
761 Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data at
762 scale. *arXiv preprint arXiv:2406.17557*, 2024.
- 763
764 Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe.
765 Lifting the curse of multilinguality by pre-training modular transformers. In Marine Carpuat, Marie-
766 Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), *Proceedings of the 2022 Conference of*
767 *the North American Chapter of the Association for Computational Linguistics: Human Language*
768 *Technologies*, pp. 3479–3495, Seattle, United States, July 2022. Association for Computational
769 Linguistics. doi: 10.18653/v1/2022.naacl-main.255. URL <https://aclanthology.org/2022.naacl-main.255>.
- 770
771 Irene Plaza, Nina Melero, Cristina del Pozo, Javier Conde, Pedro Reviriego, Marina Mayor-Rocher,
772 and María Grandury. Spanish and llm benchmarks: is mmlu lost in translation? *arXiv preprint*
773 *arXiv:2406.17789*, 2024.
- 774
775 Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Ko-
776 rhonen. Xcopa: A multilingual dataset for causal commonsense reasoning. *arXiv preprint*
777 *arXiv:2005.00333*, 2020.
- 778
779 Luiza Pozzobon, Patrick Lewis, Sara Hooker, and Beyza Ermis. From one to many: Expanding
780 the scope of toxicity mitigation in language models, 2024. URL <https://arxiv.org/abs/2403.03893>.
- 781
782 Aida Ramezani and Yang Xu. Knowledge of cultural moral norms in large language models. *arXiv*
783 *preprint arXiv:2306.01857*, 2023.
- 784
785 Abhinav Rao, Akhila Yerukola, Vishwa Shah, Katharina Reinecke, and Maarten Sap. Normad:
786 A benchmark for measuring the cultural adaptability of large language models, 2024. URL
787 <https://arxiv.org/abs/2404.12464>.
- 788
789 Yangjun Ruan, Chris J Maddison, and Tatsunori Hashimoto. Observational scaling laws and the
790 predictability of language model performance. *arXiv preprint arXiv:2405.10938*, 2024.
- 791
792 Sebastian Ruder, Noah Constant, Jan Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu,
793 Junjie Hu, Dan Garrette, Graham Neubig, et al. Xtreme-r: Towards more challenging and nuanced
794 multilingual evaluation. *arXiv preprint arXiv:2104.07412*, 2021.
- 795
796 Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. Gender bias in
797 machine translation. *Transactions of the Association for Computational Linguistics*, 9:845–874,
798 2021.
- 799
800 Lütfi Kerem Senel, Benedikt Ebing, Konul Baghirova, Hinrich Schuetze, and Goran Glavaš. Kardeş-
801 NLU: Transfer to low-resource languages with the help of a high-resource cousin – a benchmark and
802 evaluation for Turkic languages. In Yvette Graham and Matthew Purver (eds.), *Proceedings of the*
803 *18th Conference of the European Chapter of the Association for Computational Linguistics (Volume*
804 *1: Long Papers)*, pp. 1672–1688, St. Julian’s, Malta, March 2024. Association for Computational
805 Linguistics. URL <https://aclanthology.org/2024.eacl-long.100>.
- 806
807 Weiyang Shi, Ryan Li, Yutong Zhang, Caleb Ziems, Chunhua yu, Raya Horesh, Rogério Abreu
808 de Paula, and Diyi Yang. Culturebank: An online community-driven knowledge base towards cul-
809 turally aware language technologies, 2024. URL <https://arxiv.org/abs/2404.15238>.
- 810
811 Shivalika Singh, Freddie Vargus, Daniel D’souza, Börje Karlsson, Abinaya Mahendiran, Wei-Yin
812 Ko, Herumb Shandilya, Jay Patel, Deividas Mataciunas, Laura O’Mahony, Mike Zhang, Ramith
813 Hettiarachchi, Joseph Wilson, Marina Machado, Luisa Moura, Dominik Krzemiński, Hakimeh
814 Fadaei, Irem Ergun, Ifeoma Okoh, Aisha Alaagib, Oshan Mudannayake, Zaid Alyafeai, Vu Chien,
815 Sebastian Ruder, Surya Guthikonda, Emad Alghamdi, Sebastian Gehrmann, Niklas Muennighoff,
816 Max Bartolo, Julia Kreutzer, Ahmet Üstün, Marzieh Fadaee, and Sara Hooker. Aya dataset:

- 810 An open-access collection for multilingual instruction tuning. In Lun-Wei Ku, Andre Martins,
811 and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for*
812 *Computational Linguistics (Volume 1: Long Papers)*, pp. 11521–11567, Bangkok, Thailand,
813 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.620.
814 URL <https://aclanthology.org/2024.acl-long.620>.
- 815
816 Zayne Sprague, Fangcong Yin, Juan Diego Rodriguez, Dongwei Jiang, Manya Wadhwa, Prasann
817 Singhal, Xinyu Zhao, Xi Ye, Kyle Mahowald, and Greg Durrett. To cot or not to cot? chain-
818 of-thought helps mainly on math and symbolic reasoning. *arXiv preprint arXiv:2409.12183*,
819 2024.
- 820 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam
821 Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the
822 imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint*
823 *arXiv:2206.04615*, 2022.
- 824
825 Yan Tao, Olga Viberg, Ryan S Baker, and René F Kizilcec. Cultural bias and cultural alignment of
826 large language models. *PNAS nexus*, 3(9):pgae346, 2024.
- 827
828 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,
829 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models
830 based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 831
832 Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude,
833 Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargus, Phil Blunsom, Shayne
834 Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. Aya model: An
835 instruction finetuned open-access multilingual language model. In Lun-Wei Ku, Andre Martins,
836 and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for*
837 *Computational Linguistics (Volume 1: Long Papers)*, pp. 15894–15939, Bangkok, Thailand,
838 August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.845.
839 URL <https://aclanthology.org/2024.acl-long.845>.
- 840
841 Eva Vanmassenhove, Dimitar Shterionov, and Matthew Gwilliam. Machine translationese: Effects of
842 algorithmic bias on linguistic complexity in machine translation. *arXiv preprint arXiv:2102.00287*,
843 2021.
- 844
845 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue:
846 A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint*
847 *arXiv:1804.07461*, 2018.
- 848
849 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer
850 Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language
851 understanding systems. *Advances in neural information processing systems*, 32, 2019.
- 852
853 Xidong Wang, Nuo Chen, Junying Chen, Yan Hu, Yidong Wang, Xiangbo Wu, Anningzhe Gao,
854 Xiang Wan, Haizhou Li, and Benyou Wang. Apollo: An lightweight multilingual medical llm
855 towards democratizing medical ai to 6b people. *ArXiv*, abs/2403.03640, 2024a. URL <https://api.semanticscholar.org/CorpusID:268253217>.
- 856
857 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
858 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging
859 multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024b.
- 860
861 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le,
862 and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In
863 *Neural Information Processing Systems*, 2022. URL <https://api.semanticscholar.org/CorpusID:246411621>.
- 864
865 Orgest Xhelili, Yihong Liu, and Hinrich Schütze. Breaking the script barrier in multilingual
866 pre-trained language models with transliteration-based post-training alignment. *arXiv preprint*
arXiv:2406.19759, 2024.

864 L Xue. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint*
865 *arXiv:2010.11934*, 2020.
866

867 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
868 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
869 *arXiv:2407.10671*, 2024.

870 Arda Yüksel, Abdullatif Köksal, Lütfi Kerem Şenel, Anna Korhonen, and Hinrich Schütze. Turk-
871 ishmmmlu: Measuring massive multitask language understanding in turkish. *arXiv preprint*
872 *arXiv:2407.12402*, 2024.
873

874 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
875 really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.

876 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav
877 Raja, Dylan Slack, Qin Lyu, et al. A careful examination of large language model performance on
878 grade school arithmetic. *arXiv preprint arXiv:2405.00332*, 2024a.
879

880 Miaoran Zhang, Vagrant Gautam, Mingyang Wang, Jesujoba O Alabi, Xiaoyu Shen, Dietrich Klakow,
881 and Marius Mosbach. The impact of demonstrations on multilingual in-context learning: A
882 multidimensional analysis. *arXiv preprint arXiv:2402.12976*, 2024b.
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

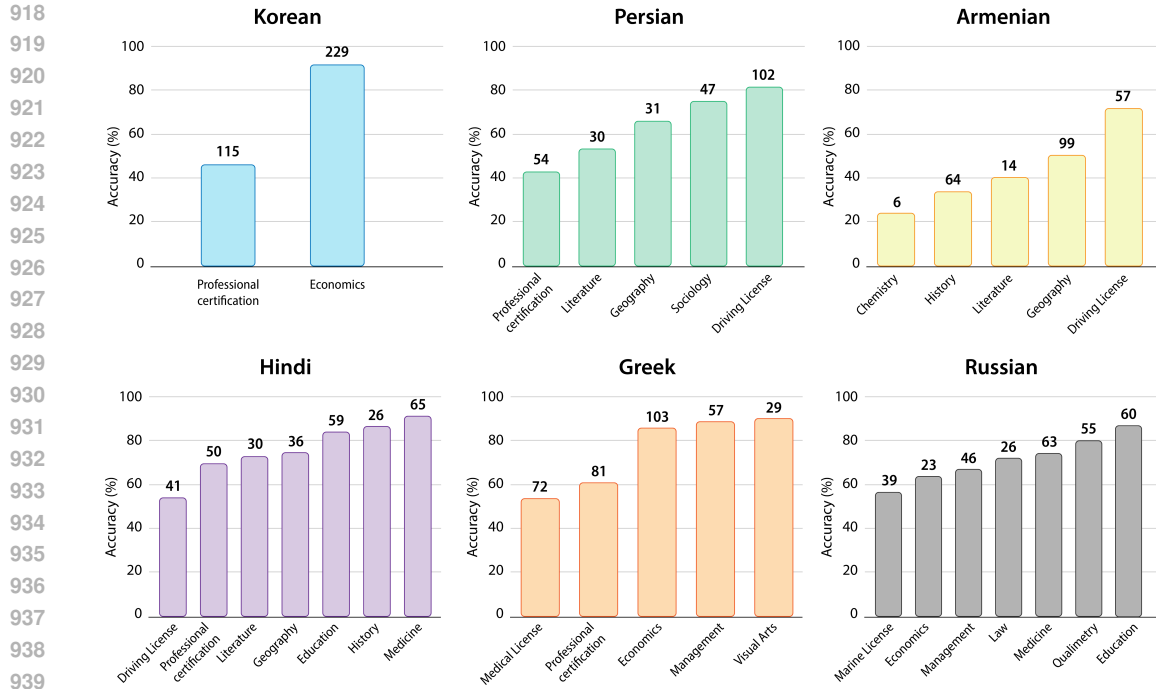


Figure 5: **GPT-4o performance across academic disciplines for Korean, Persian, Armenian, Hindi, Greek, and Russian.** Each bar is annotated with the number of questions with correct answers.

A APPENDIX

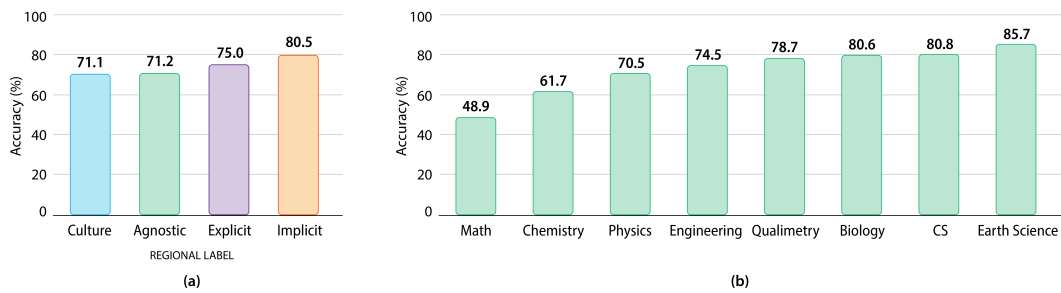


Figure 6: **GPT-4o model performance on INCLUDE-LITE.** (a) Performance across regional labels. While models typically perform better across *region-explicit* and *regional-implicit* questions, it is difficult to disentangle the difficulty of questions due to regionality from the subject matter itself (*i.e.*, *region-agnostic* questions may contain more STEM subjects that are traditionally harder for LLMs). (b) Performance across academic disciplines within STEM area. We observe models perform particularly poorly on Math and Chemistry questions.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Model	Accuracy (↑)					
	Humanities	STEM	Domain-Specific	Professional	Licenses	Average
# samples	13294	2478	1940	3189	1736	-
GPT-4o						
- 5-shot	79.0	74.2	76.8	70.1	82.1	77.5
- Zero-shot CoT	79.9	78.6	80.4	73.8	81.1	79.0
C4AI-Aya-23-35b						
- 5-shot	54.8	52.3	48.1	56.9	48.8	51.8
- Zero-shot CoT	46.0	39.1	45.9	40.6	48.2	45.7

Aya-23-8B	34.0	34.3	33.2	40.6	29.9	32.5
Mistral-7B (v0.3)	44.2	43.4	43.9	38.6	44.3	43.3
Mistral-7B-Instruct (v0.3)	44.5	42.7	43.2	40.1	43.7	43.5
Gemma-7B	55.1	53.6	55.5	47.7	62.2	54.5
Gemma-7B-Instruct	38.6	37.7	42.0	34.5	44.9	38.7
Qwen2.5-7B	53.4	54.2	59.1	51.3	57.8	54.1
Qwen2.5-7B-Instruct	53.5	53.3	58.1	49.5	58.6	53.7
Qwen2.5-14B	61.4	60.9	66.0	57.1	65.1	61.4
Llama-3-8B	51.7	49.8	52.1	43.4	51.3	50.3
Llama-3-8B-Instruct	50.7	46.9	52.9	44.3	54.4	49.8
XGLM-7.5B	24.6	26.8	24.9	25.1	25.7	25.0
BLOOM-7.1B	25.1	27.2	26.0	25.1	24.0	25.3
BLOOMZ-7.1B	26.6	30.0	25.7	25.3	30.6	27.0

Table 2: Accuracy performance of GPT-4o on **INCLUDE-LITE** grouped by high-level topics. Where **Humanities** include Social Science, Humanities, and General knowledge. **STEM** includes Applied Science and STEM. **Domain-specific** covers Business & Commerce and Health oriented education. **Professional** includes professional certifications. **Licenses** cover Marine, Fishing, and Driving licenses.

Language	Academic Humanities	Academic STEM studies	Academic Domain-specific studies	Professional	License	Avg (%)
Albanian	95.0	88.0	83.5	-	-	89.50
Arabic	77.8	82.0	80.5	-	76.2	78.30
Armenian	52.7	32.0	-	-	72.2	53.60
Azerbaijani	71.3	73.6	71.4	-	-	71.90
Basque	-	-	-	64.8	-	64.80
Belarusian	51.8	42.0	-	-	-	50.90
Bengali	71.1	90.0	-	84.3	-	76.80
Bulgarian	93.8	60.0	-	-	-	90.70
Chinese	71.5	66.7	58.2	52.1	84.5	66.10
Croatian	89.0	82.0	-	-	-	88.40
Dutch; Flemish	86.6	87.5	80.0	-	-	86.40
Estonian	90.7	98.0	100.0	-	-	92.40
Finnish	67.0	87.0	77.8	-	-	69.90
French	83.8	50.0	81.2	-	68.1	80.70
Georgian	87.6	-	-	-	-	87.60
German	62.6	64.0	-	-	87.0	66.90
Greek	84.7	84.0	89.2	58.6	-	71.50
Hebrew	62.0	-	-	-	88.6	86.20
Hindi	77.7	71.9	91.5	71.8	57.7	75.10
Hungarian	66.3	80.6	-	-	-	75.80
Indonesian	84.0	69.1	-	84.8	-	79.50
Italian	87.7	87.2	91.7	95.5	-	90.00
Japanese	-	-	-	78.1	96.0	81.60
Kazakh	80.4	-	-	-	-	80.40
Korean	91.6	-	-	46.4	-	69.00
Lithuanian	92.0	97.1	82.5	81.2	-	90.60
Malay	84.5	-	80.3	-	-	83.00
Malayalam	69.6	66.0	55.0	-	80.9	70.80
Nepali	-	-	-	61.6	83.2	72.40
Macedonian	96.0	86.0	89.3	-	-	92.40
Persian	66.0	25.0	-	49.6	81.6	64.60
Polish	100.0	64.6	-	80.0	-	78.80
Portuguese	84.7	63.3	67.9	-	-	76.40
Serbian	92.2	86.0	-	-	-	91.60
Spanish	83.6	88.0	96.0	-	-	84.40
Tagalog	86.8	-	-	-	90.7	87.40
Tamil	70.6	54.0	-	-	-	69.10
Telugu	66.9	70.7	-	-	-	68.20
Turkish	62.0	52.0	75.9	-	-	65.30
Ukrainian	85.8	84.0	-	-	-	85.60
Urdu	61.7	65.3	100.0	-	-	62.50
Uzbek	63.6	84.0	-	73.3	-	69.70
Vietnamese	84.4	86.0	-	-	-	84.50
Russian	77.5	83.4	70.8	-	63.9	75.00

Table 3: Accuracy performance of GPT-4o (5-shot) on **INCLUDE-LITE** for each language. **Hu-**manities include Social Science, Humanities, and General knowledge. **STEM** includes Applied Science and STEM. **Domain-specific** covers Business & Commerce and Health oriented education. **Professional** includes professional certifications. **Licenses** cover Marine, Fishing, and Driving licenses.

A.1 COLLECTED LANGUAGES

Table 4 provides information about the languages in **INCLUDE**.

Language	Script	Family	Branch	Availability	Count
Albanian	latin	Indo-European	Albanian	Mid	2365
Amharic	ge'ez	Afro-Asiatic	Semitic	Low	131
Arabic	perso-arabic	Afro-Asiatic	Semitic	High	15137
Armenian	armenian	Indo-European	Armenian	Low	1669
Assamese	bengali-assamese	Indo-European	Indo-Iranian	Low	323
Azerbaijani	latin	Turkic	Azerbaijani North	Mid	6937
Basque	latin	Isolate		Low	719
Belarusian	cyrillic	Indo-European	Slavic East	Low	687
Bengali	bengali-assamese	Indo-European	Indo-Iranian	Mid	15259
Bulgarian	cyrillic	Indo-European	Slavic South Eastern	Mid	2937
Chinese	chinese	Sino-Tibetan	Chinese	High	12977
Croatian	latin	Indo-European	Slavic South Western	Mid	2879
Czech	latin	Indo-European	Slavic West	High	50
Danish	latin	Indo-European	Germanic	Mid	732
Dutch; Flemish	latin	Indo-European	Germanic	High	2222
Estonian	latin	Uralic	Finnic	Mid	952
Finnish	latin	Uralic	Finnic	Mid	1574
French	latin	Indo-European	Italic	High	2457
Georgian	mkherduli	Kartvelian	Georgian	Low	599
German	latin	Indo-European	Germanic	High	1590
Greek	greek	Indo-European	Greek	Mid	6570
Hebrew	hebrew	Afro-Asiatic	Semitic	Mid	2457
Hindi	devanagari	Indo-European	Indo-Iranian	Mid	5167
Hungarian	latin	Uralic	Hungarian	Mid	2267
Indonesian	latin	Austronesian	Malayo-Polynesian	High	12013
Italian	latin	Indo-European	Italic	High	3038
Japanese	kanji	Japonic	Japanese	High	2699
Kannada	kannada	Dravidian	Southern	Low	335
Kazakh	cyrillic	Turkic	Western	Low	5736
Korean	hangul	Koreanic	Korean	Mid	1781
Lithuanian	latin	Indo-European	Eastern Baltic	Mid	1397
Malay	latin	Austronesian	Malayo-Polynesian	Mid	1021
Malayalam	vatteluttu	Dravidian	Southern	Low	275
Marathi	devanagari	Indo-European	Indo-Iranian	Mid	313
Nepali	devanagari	Indo-European	Indo-Iranian	Mid	1470
Macedonian	cyrillic	Indo-European	Slavic South Eastern	Low	2075
Oriya	odia	Indo-European	Indo-Iranian	Low	241
Punjabi; Punjabi	gurmukhi	Indo-European	Indo-Iranian	Low	453
Persian	perso-arabic	Indo-European	Indo-Iranian	High	23990
Polish	latin	Indo-European	Slavic West	High	2023
Portuguese	latin	Indo-European	Italic	High	1407
Russian	cyrillic	Indo-European	Slavic East	High	10169
Serbian	cyrillic	Indo-European	Slavic South	Mid	1636
Sinhala; Sinhalese	sinhala	Indo-European	Indo-Iranian	Low	325
Slovak	latin	Indo-European	Slavic West	Mid	131
Spanish	latin	Indo-European	Italic	High	2559
Swedish	latin	Indo-European	Germanic	Mid	5102
Tagalog	latin	Austronesian	Malayo-Polynesian	Low	530
Tamil	tamil	Dravidian	Southern	Mid	945
Telugu	telugu	Dravidian	South-Central	Low	11568
Turkish	latin	Turkic	Southern	High	2710
Ukrainian	cyrillic	Indo-European	Slavic East	Mid	1482
Urdu	perso-arabic	Indo-European	Indo-Iranian	Low	122
Uzbek	latin	Turkic	Eastern	Low	2878
Vietnamese	latin	Austro-Asiatic	Mon-Khmer	High	8901

Table 4: Languages in **INCLUDE** with their associated metadata and the total count of the samples per language.

1134 A.2 DETAILS ON EXAM SOURCES PARSING

1135
1136 Figure 9 presents the questionnaire we distributed to the community to gather a diverse set of
1137 multiple-choice exams. It was distributed among university student organizations and researchers at
1138 our institution.⁶ Participation was voluntary and not incentivized.

1140 A.3 PERFORMANCE ACROSS ACADEMIC AREAS AND FIELDS

1141
1142 Distribution of academic areas and academic fields with the respective number of questions is
1143 presented in Figure 7. GPT-4o performance across languages and academic areas is in Table 5.
1144 GPT-4o performance across languages, academic fields, and related regional features is in Tables 6
1145 and 7.

1147 A.4 REGIONAL LABELS: ANNOTATION

1148
1149 First, we categorized the exams into one of eight broad academic areas, *e.g.*, Humanities or Social
1150 Sciences, and then further classified each exam into a specific academic fields, *e.g.*, History or
1151 Geography. This categorization was done manually, taking into account both the exam’s learning
1152 level and the exam’s original topic.

1153 Building on these categories, we applied one of four labels—agnostic, culture-related, region-explicit,
1154 or region-implicit—based on the degree of dependence on localized knowledge required to answer
1155 the exam questions. The labels reflect the extent to which specific cultural or regional knowledge is
1156 necessary. Table 8 provides examples illustrating how different exams were typically classified under
1157 each label, to show the relationship between categories and labels.

1158 The “region implicit” label was applied when we suspected that exam content might vary across
1159 regions but could not reliably detect specific regional differences. For example, historical events,
1160 literary works, and religious interpretations may differ significantly depending on the region. Similarly,
1161 fields like marketing, management, social work, and insurance—though rooted in shared theoretical
1162 foundations—can be practiced differently across regions. When we encountered such uncertainty, we
1163 labeled the subject as “region implicit.”

1164 Within the Humanities, fields such as Visual Arts, History, Philosophy, Religious Studies, Performing
1165 Arts, Culturology, and Literature were labeled “region implicit” when the content was not explicitly
1166 tied to a particular region. However, if the exam was region-specific (*e.g.*, Greek literature), we
1167 categorized it as “culture-related.”

1168 In the Social Sciences, Psychology was classified as “region explicit” if the exam focused on regional
1169 clinical practices; otherwise, it was “region implicit” when dealing with broader psychological
1170 theories that may vary across regions. Geography was labeled “region implicit” if the exam involved
1171 political geography and “agnostic” if it focused on general geographic knowledge. Similarly, disci-
1172 plines like Sociology, Political Science, and Anthropology were classified as either “region implicit”
1173 or “culture-related,” depending on whether regional specificity was required, much like History.

1174 For Economics, exams were labeled “agnostic” when covering general economic theories, “region
1175 explicit” when addressing regional regulations, and “region implicit” when regional applications
1176 were uncertain. In STEM fields, most disciplines were categorized as “agnostic,” with the exception
1177 of Qualimetry, which was labeled “region explicit” due to its specific application in post-Soviet
1178 countries for quantitative and qualitative assessment according to regional standards.⁷

1179 Exams related to theoretical medical subjects, such as Anatomy, were classified as “agnostic.” In
1180 contrast, exams covering clinical practices and guidelines specific to a region were labeled as “region
1181 explicit,” while others were marked as “region implicit” if regional dependence was unclear.

1182 Accounting is generally tied to region-specific practices, so it was consistently classified as “region
1183 explicit.” Other disciplines within the Business and Commerce category were treated similarly to
1184 Economics and labeled as mostly as “region implicit.” In some cases, there were “agnostic” and
1185 “region explicit” exams.

1186
1187 ⁶We will provide institutional details upon paper publication.

⁷<https://en.wikipedia.org/wiki/Qualimetry>

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
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

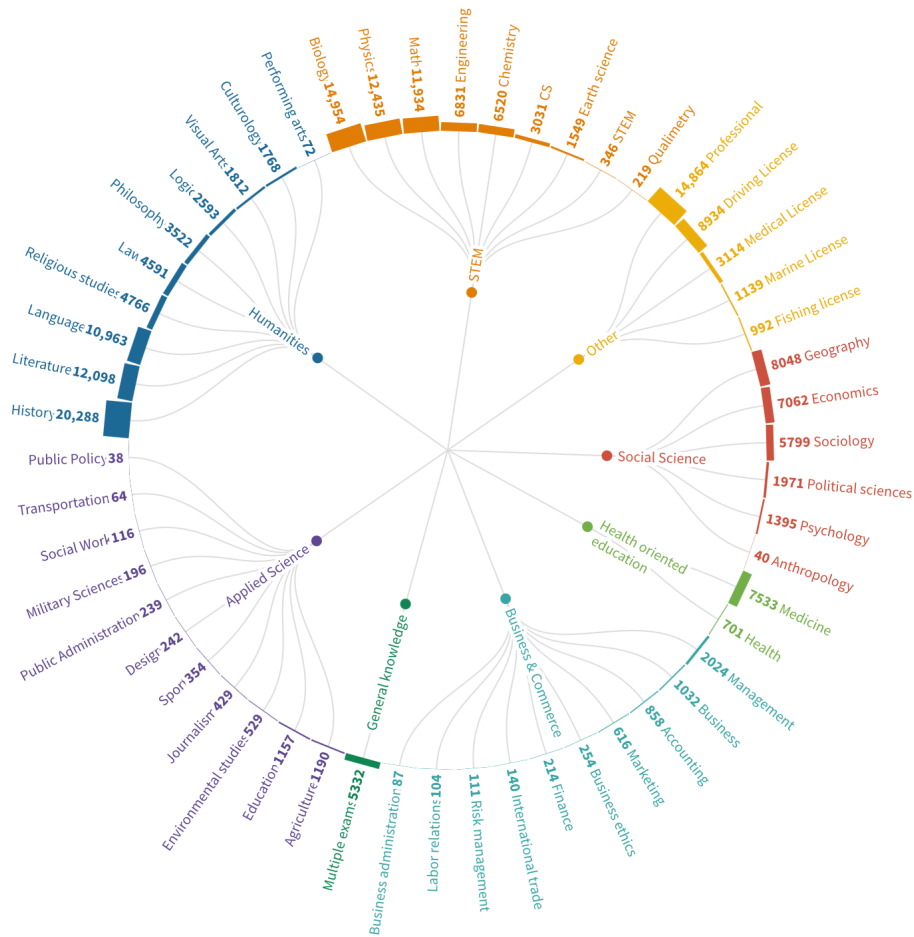


Figure 7: Academic domain and academic fields with the number of examples across all languages.

Exams in Applied Science disciplines were typically categorized as “region implicit” due to the potential involvement of regional variations. Similarly, exams in Military Sciences, Public Administration, and Public Policy were marked as “region explicit” when tied to specific regions (*e.g.*, Basics of National Security of the Republic of Azerbaijan) and “region implicit” when regional specifics were less pronounced. For exams focused on theoretical aspects, we used the “agnostic” label (*e.g.*, Theoretical Foundations of Food Engineering in Agriculture).

Finally, exams covering multiple topics were classified as “region implicit” unless they explicitly focused on cultural aspects of a particular region, in which case they were labeled as “culture-related.”

A.5 IMPLEMENTATION DETAILS

Each model was evaluated using a single A100 GPU (80GB memory), with evaluation times averaging approximately 4 hours for **INCLUDE-LITE**. For all models, we set the decoding temperature to 0, prioritizing deterministic outputs.

We configured response context windows based on model size and task requirements. For models such as Aya-23-8B, Mistral-7B (v0.3), Mistral-7B-Instruct (v0.3), Gemma-7B, Gemma-7B-Instruct, Qwen2.5-7B, Qwen2.5-7B-Instruct, Llama-3-8B, Llama-3-8B-Instruct, XGLM-7.5B, BLOOM-7.1B,

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

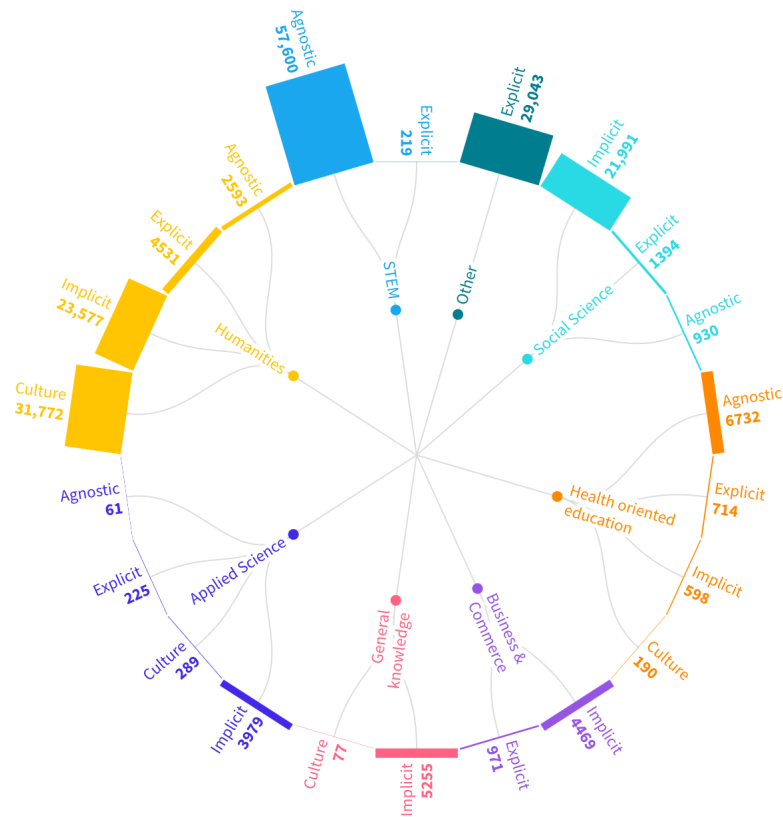


Figure 8: **Distribution of regional labels per Academic area.**

and BLOOMZ-7.1B, we set a window size of 40 tokens. Larger models, including C4AI-Aya-23-35B and GPT-4, were evaluated using a 512-token context window for 5-shot tasks and 1024 tokens for zero-shot chain-of-thought (CoT) reasoning.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Collection of Multilingual Exams 🌍👉

Thank you for participating in this survey!

We invite you to help us build the next era of multilingual Large Language Models (LLMs)!

We are collecting a pool of **multilingual exams** from various languages and cultural sources. Multilingual data is crucial in AI to ensure inclusivity and fairness, enabling systems to understand and serve diverse cultures and communities equally.

In this Google Form, we ask you to provide a list of examinations in your native language.

These sets of **question-answer pairs** can fall under the categories below:

- **Educational exams:** National examinations of your country or exams relative to the level of education: e.g., High school Physics, College/University Literature.
- **Professional exams:** e.g., Law - Bar examination, Medical License examinations.
- **Practical tests:** e.g., Driver license test, Marine License Practice Tests.
- Any other examination that tests general knowledge

Please provide the following information per provided examination

- **Name:** The name of the exam (e.g., "Physics exam - High school")
- **Description:** A short description of the examination
- **Language** of the data
- **URL** of the data: The link where the data is available (if applicable).
- **Answers:** Describe how the answers are provided, i.e. inside the question file, as a separate file (provide the url of this file) etc.
- Additional information, if available, such as the approximate size (in number of rows), a short description of the nature of the data and their format, etc.

Here are some indicative examples:

- **Name:** Greek national examination on different subjects
Description: PDFs of published previous exams for all academic topics.
Language: Greek
URL of the source: <https://www.panellines.net/πανελληνιες-εξετασεις/>
Answers: In a separate file under the same URL
Size: ~ 45 academic topics, for the last 20 years

- **Name:** Iranian Law Examination:
Description: The Bar examination in Iran.
Language: Farsi
URL: https://www.hevyalaw.com/uploads/files/1934640179vokala_hoghoqhi_1402_.pdf
Answers: At the end of each question file

🎓 Educational examinations
 National examinations of your country or exams relative to the level of education.

Please provide a list of examinations based on the format above (i.e. name, description, language, url etc.)

Long-answer text

👨‍⚖️ Professional examinations
 e.g., Law - Bar examination, Medical License examinations, etc.

Please provide a list of examinations based on the format above (i.e. name, description, language, url etc.)

Long-answer text

🛠️ Practical examinations
 e.g., Driver license test, Marine License Practice Tests, etc.

Figure 9: Exam source collection form sent to the academic community.

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Language	Academic Area	Accuracy	Count
Albanian	Humanities	95.1	223
	Business & Commerce	85.7	223
	Social Science	94.5	55
Arabic	Humanities	79.0	105
	Business & Commerce	79.3	82
	General Knowledge	86.7	105
	Other	76.2	105
	STEM	82.0	50
Armenian	Social Science	67.6	105
	Humanities	34.7	225
	Other	72.2	79
Azerbaijani	STEM	28.0	50
	Social Science	50.5	196
	Applied Science	75.9	108
Basque	Humanities	74.1	108
	Business & Commerce	62.5	96
	Health-Oriented Education	80.2	96
	Social Science	67.6	108
Belarusian	Other	64.8	500
	Humanities	50.8	490
Bengali	STEM	42.0	50
	Humanities	62.0	166
	General Knowledge	80.1	166
	Other	84.3	166
Bulgarian	STEM	88.0	50
	Humanities	96.4	250
	Social Science	60.0	50
Chinese	Social Science	91.2	250
	Applied Science	73.2	71
	Humanities	67.8	87
	Business & Commerce	53.5	71
	Health-Oriented Education	60.9	87
Croatian	Other	68.3	142
	Social Science	76.1	71
	Humanities	86.8	250
	STEM	82.0	50
Dutch; Flemish	Social Science	90.8	250
	Humanities	86.0	243
Estonian	Social Science	86.8	243
	Humanities	90.1	161
Finnish	STEM	97.2	36
	Humanities	69.5	226
	Health-Oriented Education	75.6	45
French	Social Science	64.6	226
	Humanities	86.5	266
	Other	68.1	47
Georgian	Social Science	74.3	74
	Humanities	87.6	500
German	Social Science	62.6	91
	Humanities	83.8	37
	Business & Commerce	89.1	64
	Other	57.5	266
Greek	Social Science	84.2	133
	Humanities	60.0	50
	Other	88.6	500
Hebrew	Applied Science	83.1	71
	Humanities	72.9	96
	General Knowledge	83.1	71
	Health-Oriented Education	91.5	71
	Other	64.1	142
Hindi	Social Science	74.6	71
	Applied Science	79.8	341
	Social Science	66.3	184
Hungarian	Applied Science	71.2	125
	Humanities	82.4	125
	Other	83.2	125
	STEM	60.0	50
	Social Science	84.8	125
Indonesian	Applied Science	85.7	35
	Humanities	85.0	167
	Other	95.5	155
	Social Science	89.8	167

Language	Academic Area	Accuracy	Count
Japanese	Other	80.2	501
	Humanities	80.4	500
Kazakh	Other	46.0	250
	Social Science	91.6	250
Korean	Humanities	91.6	335
	Business & Commerce	77.5	40
	Other	81.2	48
Lithuanian	STEM	97.1	34
	Social Science	93.5	77
	Humanities	84.3	178
Malay	Business & Commerce	79.8	178
	Social Science	84.8	145
Malayalam	Humanities	64.3	56
	General Knowledge	73.1	78
	Health-Oriented Education	55.0	100
	Other	80.9	194
Nepali	STEM	66.0	47
	Other	72.4	500
Macedonian	Humanities	96.9	224
	Business & Commerce	89.3	224
	STEM	86.0	50
	Social Science	92.5	53
Persian	Humanities	55.3	141
	Other	62.4	250
	Social Science	74.5	141
Polish	Other	80.0	496
	STEM	62.5	48
Portuguese	Applied Science	58.3	84
	Humanities	81.8	154
	Business & Commerce	56.9	84
	Health-Oriented Education	67.1	67
	Other	67.6	169
Russian	Applied Science	87.0	69
	Humanities	76.8	69
	Business & Commerce	66.7	69
	Health oriented education	74.1	85
	Other	63.9	97
	STEM	80.9	94
Serbian	Social Science	76.8	69
	Humanities	90.4	313
	STEM	84.0	50
Spanish	Social Science	95.2	187
	Humanities	77.2	250
	Health oriented education	96.0	25
Tagalog	STEM	88.0	25
	Social Science	89.6	250
	Humanities	86.8	425
Tamil	Other	90.7	75
	General knowledge	70.6	500
Telugu	STEM	54.0	50
	Applied Science	73.5	166
	Humanities	66.0	191
Turkish	Social Science	66.9	166
	Humanities	62.0	166
	Business & Commerce	75.9	166
Ukrainian	STEM	52.0	50
	Social Science	62.0	166
	Humanities	92.4	250
Urdu	STEM	84.0	50
	Social Science	79.2	250
Uzbek	Humanities	61.7	300
	STEM	63.3	49
	Humanities	62.9	240
Vietnamese	Other	73.3	240
	STEM	84.0	50
	Social Science	71.4	21
Vietnamese	Humanities	88.0	250
	STEM	86.0	50
	Social Science	80.8	250

Table 5: GPT-4o (5-shot) performance on INCLUDE-LITE per language and academic area. Areas with less than 30 examples were excluded from the analysis.

Language	Academic Field	Regional Feature	Accuracy	Count
1404				
1405				
1406				
1407				
1408				
1409				
1410				
1411				
1412				
1413				
1414				
1415				
1416				
1417				
1418				
1419				
1420				
1421				
1422				
1423				
1424				
1425				
1426				
1427				
1428				
1429				
1430				
1431				
1432				
1433				
1434				
1435				
1436				
1437				
1438				
1439				
1440				
1441				
1442				
1443				
1444				
1445				
1446				
1447				
1448				
1449				
1450				
1451				
1452				
1453				
1454				
1455				
1456				
1457				

Table 6: GPT-4o (5-shot) performance on **INCLUDE-LITE** per language, academic field, and regional label. Fields with less than 30 examples were excluded from the analysis (Part 1)

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Language	Academic Field	Regional Feature	Accuracy	Count
Hungarian	Agriculture	Implicit	82.4	170
	Architecture and Design	Explicit	85.7	42
	Environmental Studies and Forestry	Implicit	74.4	129
	Economics	Implicit	80.8	78
Indonesian	Geography	Implicit	48.1	81
	Human Physical Performance and Recreation	Implicit	71.2	125
	Language	Culture	79.5	78
	Professional Certification	Region explicit	83.2	125
	Economics	Region explicit	77.8	36
Italian	Geography	Implicit	87.5	32
	Sociology	Implicit	87.7	57
	Agriculture	Implicit	85.7	35
	History	Implicit	90.4	94
Japanese	Professional Certification	Region explicit	95.5	155
	Psychology	Implicit	95.0	60
	Sociology	Implicit	87.7	65
Kazakh	Driving License	Region explicit	96.0	99
	Medical License	Region explicit	86.1	201
	Professional Certification	Region explicit	66.7	201
Korean	History	Culture	78.4	241
	History	Implicit	94.9	79
	Literature	Culture	76.7	180
Lithuanian	Professional Certification	Region explicit	46.0	250
	Economics	Implicit	91.6	250
Malay	History	Implicit	91.6	335
	Finance	Implicit	77.5	40
	Professional Certification	Region explicit	81.2	48
	Earth Science	Agnostic	97.1	34
Malayalam	Economics	Implicit	93.5	77
	History	Implicit	84.3	178
	Accounting	Region explicit	79.8	178
Nepali	Geography	Implicit	85.3	129
	History	Implicit	61.5	52
	Multiple Exams	Culture	72.7	77
	Health	Implicit	55.0	100
North Macedonian	Marine License	Explicit	80.9	194
	Driving License	Explicit	83.2	250
	Professional Certification	Explicit	61.6	250
Persian	History	Implicit	95.8	48
	Philosophy	Implicit	97.3	74
	Visual Arts	Implicit	97.1	102
	Business	Implicit	89.3	224
	Sociology	Implicit	92.5	53
Polish	Literature	Culture	51.6	31
	Driving License	Explicit	81.6	125
	Professional Certification	Explicit	43.2	125
	Geography	Implicit	66.0	47
Portuguese	Sociology	Implicit	74.6	63
	Professional Certification	Explicit	80.0	496
	Math	Agnostic	61.7	47
Russian	Agriculture	Implicit	70.0	40
	Philosophy	Implicit	83.3	84
	Management	Implicit	57.9	57
	Health	Implicit	70.3	37
	Economics	Implicit	89.7	126
Serbian	Education	Implicit	87.0	69
	Law	Explicit	72.2	36
	Management	Implicit	66.2	65
	Medicine	Explicit	73.3	60
	Marine License	Explicit	56.5	69
	Qualimetry	Explicit	79.7	69
Spanish	Economics	Implicit	63.9	36
	History	Implicit	91.5	235
	Philosophy	Implicit	87.5	56
	Psychology	Implicit	99.2	125
Tagalog	Sociology	Implicit	91.1	45
	Language	Culture	69.6	46
	Law	Explicit	67.0	109
	Literature	Implicit	93.8	64
	Philosophy	Implicit	90.3	31
Tamil	Economics	Explicit	95.6	91
	Geography	Implicit	86.2	159
	Culturology	Culture	91.6	203
	History	Culture	85.3	116
Telugu	Language	Culture	79.2	106
	Driving License	Explicit	90.7	75
	Multiple Exams	Implicit	70.6	500
	Education	Implicit	73.0	100
	History	Culture	64.7	119
Turkish	History	Implicit	63.9	36
	Economics	Explicit	60.0	45
	Geography	Implicit	73.2	82
	Political Sciences	Implicit	63.3	30
	History	Implicit	71.2	73
Ukrainian	Philosophy	Implicit	74.6	63
	Business	Implicit	75.9	166
	Geography	Implicit	53.8	130
	Sociology	Implicit	91.7	36
Urdu	Law	Explicit	92.4	250
	Physics	Agnostic	84.0	50
	Psychology	Implicit	79.2	250
Uzbek	Culturology	Culture	61.7	300
	History	Implicit	66.1	124
	Law	Explicit	60.6	109
Vietnamese	Medical License	Explicit	73.3	240
	History	Implicit	88.3	239
Geography	Implicit	80.8	250	

Table 7: GPT-4o (5-shot) performance on **INCLUDE-LITE** per language, academic field, and regional label. Fields with less than 30 examples were excluded from the analysis (Part 2)

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

Academic area	Academic field	Label
Humanities	Logic	Agnostic
	Law	Region Explicit
	Language	Culture
	Visual Arts, History, Philosophy, Religious studies, Performing arts, Culturology, Literature	Region implicit/ Culture
Social Science	Sociology, Political sciences, Anthropology	Region implicit/Culture
	Economics	Region implicit/Agnostic/Region explicit
	Psychology	Region implicit/Region explicit
	Geography	Region implicit/Agnostic
STEM	Math, Physics, CS, Biology, Earth science, Chemistry, Engineering	Agnostic
	Qualimetry	Region explicit
	Medicine	Agnostic/Region implicit/Region explicit
Health oriented education	Health	Region implicit/Region explicit
	Accounting	Region explicit
Business and Commerce	Management, Marketing, Industrial and labor relations, International trade, Risk management and insurance, Business administration, Business ethics, Business, Finance	Region implicit/Region explicit/Agnostic
	Agriculture, Library and museum studies, Transportation	Region implicit/Agnostic
	Military Sciences, Public Administration, Public Policy	Region implicit/Region explicit
Applied Science	Architecture and Design, Family and consumer science, Environmental studies and forestry, Education Journalism, media studies, and communication, Social Work, Human physical performance and recreation	Region implicit
	Driving license, Marine license, Fishing license, Medical license, Public administration, Professional certification	Region explicit
	General knowledge	Multiple exams

Table 8: Annotation schema for high-level **Academic area** and fine-grained **Academic field**. The **Label** column lists the most likely *regionality* label for these exams in our dataset (e.g., *region-agnostic*, *implicit*, *explicit*) or *cultural*), though all exams from which we collect data are individually labeled with a *regionality* category. The first label is the most frequent one.

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

Language	Acc ($k:50$)	Acc ($k:100$)	Acc ($k:200$)	Acc ($k:512$)	Total gain
Uzbek	51.4	60.6	66.6	68.6	17.2
Armenian	28.0	30.7	36.0	41.1	13.1
Malayalam	57.0	57.4	61.0	69.9	12.9
Urdu	53.7	56.8	58.8	62.2	8.5
Greek	58.0	58.2	63.8	66.4	8.4
Korean	60.4	61.0	62.4	68.8	8.4
Chinese	57.2	61.8	63.5	65.5	8.3
Finnish	63.3	64.4	67.0	69.1	5.8
Basque	60.0	60.8	63.8	64.8	4.8
Polish	74.1	75.2	75.4	78.1	4.0
Azerbaijani	67.7	69.2	70.4	71.5	3.8
Dutch; Flemish	81.9	82.9	83.8	85.3	3.4
Telugu	63.9	63.9	64.8	66.6	2.7
Hindi	72.0	72.4	73.7	74.4	2.4
German	64.0	65.5	65.5	66.2	2.2
Malay	80.6	81.8	82.4	82.8	2.2
Tamil	67.3	67.3	67.8	69.5	2.2
Arabic	76.3	76.8	77.9	78.4	2.1
russian	72.6	73.6	74.1	74.6	2.0
Italian	88.0	88.5	89.2	89.6	1.6
Spanish	82.4	83.1	83.3	84.0	1.6
Japanese	78.6	78.6	79.4	80.0	1.4
Georgian	86.2	86.4	87.0	87.6	1.4
Vietnamese	82.4	82.5	84.9	83.8	1.4
Turkish	63.5	64.1	64.4	64.8	1.3
Kazakh	79.2	79.6	80.4	80.4	1.2
Portuguese	72.8	73.5	73.5	74.0	1.2
Bengali	75.2	75.4	76.1	76.3	1.1
Persian	60.9	61.1	61.3	61.9	1.0
Belarusian	49.5	50.0	50.0	50.2	0.7
French	80.0	80.2	80.4	80.7	0.7
Indonesian	77.8	78.2	78.4	78.5	0.7
Albanian	88.9	89.3	89.3	89.5	0.6
Lithuanian	89.7	89.7	90.1	90.3	0.6
Estonian	92.0	92.0	92.4	92.4	0.4
Croatian	87.8	88.0	88.2	88.0	0.2
Hungarian	75.3	75.3	75.5	75.5	0.2
Nepali	71.8	72.0	71.6	72.0	0.2
Bulgarian	90.7	90.7	90.7	90.7	0.0
Hebrew	86.0	86.0	86.0	86.0	0.0
Macedonian	92.4	92.4	92.4	92.4	0.0
Serbian	91.5	91.5	91.5	91.5	0.0
Tagalog	87.4	87.4	87.4	87.4	0.0
Ukrainian	85.5	85.5	85.5	85.5	0.0

Table 9: GPT-4o performance for different values of k (the output generation length) per language on INCLUDE-LITE and total performance gain from $k = 50$ to 512.