

# AraEval: An Arabic Multi-Task Evaluation Suite for Large Language Models

Anonymous ACL submission

## Abstract

The rapid advancements of Large Language models (LLMs) necessitate robust benchmarks. In this paper, we present AraEval, a pioneering and comprehensive evaluation suite specifically developed to assess the advanced knowledge, reasoning, truthfulness, and instruction following capabilities of foundation models within the Arabic context. AraEval includes a diverse set of evaluation tasks that test various dimensions of knowledge and reasoning, with a total of 24,378 samples. These tasks cover areas such as linguistic understanding, factual recall, logical inference, commonsense reasoning, mathematical problem-solving, and domain-specific expertise, ensuring that the evaluation goes beyond basic language comprehension. It covers multiple domains of knowledge, such as science, history, religion, and literature, ensuring that the LLMs are tested on a broad spectrum of topics relevant to Arabic-speaking contexts. AraEval is designed to facilitate comparisons across different foundation models, enabling LLM developers and users to benchmark performance effectively. In addition, it provides diagnostic insights to identify specific areas where models excel or struggle, guiding further development. Datasets and evaluation integration can be found at [\[https://redacted/for/anon/sub\]](https://redacted/for/anon/sub)

## 1 Introduction

With the unprecedented scaling of large language models (LLMs) (OpenAI, 2022; Google, 2024; Anthropic, 2022; Dubey et al., 2024; Mistral, 2024; Team et al., 2024; Liu et al., 2024; Team, 2024), algorithmic intelligence has reached new frontiers (Guo et al., 2025; Jaech et al., 2024) across numerous domains, demonstrating remarkable abilities in tasks ranging from creative writing (Gómez-Rodríguez, 2023), program synthesis (Jimenez et al., 2023; Khan et al., 2024), instruction following (Zhou et al., 2023), knowledge extraction

(Hendrycks et al., 2021; Wang et al., 2024b) to rich scientific reasoning (Mialon et al., 2023; Rein et al., 2023). The field has witnessed breakthroughs, with models matching or surpassing expert human performance (Glazer et al., 2024) - from solving olympiad-level problems (AlphaCode Team, 2023; Chervonyi et al., 2025) to generating research-level insights (Google, 2025; OpenAI, 2025) - catalyzing massive industry investments<sup>1</sup> and research efforts (Workshop et al., 2022; Lovenia et al., 2024; LAION-AI, 2025; Lozhkov et al., 2024). As model capabilities rapidly expand and emerge on a different scale (Wei et al., 2022; Srivastava et al., 2022), systematic evaluation (Laskar et al., 2023; Phan et al., 2025) serves as a vital proxy for decision making across the ecosystem, enabling key stakeholders - from developers and regulators to investors, researchers, and industry practitioners - to make informed strategic choices (Handa et al., 2025) about model development, deployment, and adoption (Latent Space, 2024).

Despite progress, the evaluation landscape remains significantly skewed towards English and other high-resource languages (Joshi et al., 2020), creating a significant gap in our understanding of LLM capabilities in different linguistic and cultural contexts. In addition to that Yong et al. (2023) showed that safety or instruction following don't generalize with low-resource languages. This disparity is particularly pronounced for Arabic, the fifth most spoken language worldwide with more than 400 million speakers (Eberhard et al., 2020) and rich dialectal variations spanning more than 20 countries. Although recent years have seen the emergence of Arabic-specific language models (Bari et al., 2024; Abbas et al., 2025; Sengupta et al., 2023b; Huang et al., 2023) and the increasing integration of Arabic in multilingual models (Team,

<sup>1</sup><https://openai.com/index/announcing-the-stargate-project/>

2024; Mistral, 2024; Jaech et al., 2024), comprehensive evaluation frameworks for assessing their capabilities remain limited.

Existing Arabic evaluation efforts have primarily focused on translating english benchmarks (Huang et al., 2023; OpenAI, 2025; Sengupta et al., 2023b) or targeted towards only knowledge base questions (Koto et al., 2024; Almazrouei et al., 2023), lacking the systematic multi-task assessment necessary for understanding model performance across diverse linguistic phenomena and real-world applications. Notable initiatives like ArabicMMLU (Koto et al., 2024), Exams (Hardalov et al., 2020), ACVA (Huang et al., 2023), Belebele (Bandarkar et al., 2023), and AraDiCE (Mousi et al., 2024), along with various leaderboard efforts (El Filali et al., 2024), have established foundational work in Arabic language evaluation. However, these benchmarks predominantly focus on language comprehension, dialectal understanding, and knowledge retrieval tasks. As LLM capabilities rapidly evolve, there is a pressing need for evaluating more emergent capabilities such as complex reasoning, generation, instruction following and sophisticated domain-specific applications (Laskar et al., 2024) for fine grained *System 2 Thinking* evaluation (Kahneman, 2011). Recent work by Bari et al. (2024) and Abbas et al. (2025) have attempted to address these limitations through human evaluation, but this approach faces inherent challenges of *scalability* and *consistency*, being vulnerable to variations in setup, prompt design, individual assessor biases, and temporal factors.

This evaluation gap poses significant challenges for the development and deployment of Arabic Language Technologies (ALT). In this work, we introduce AraEval, a comprehensive Arabic multi-task evaluation suite designed to rigorously assess large language models (LLMs) in Arabic. AraEval introduces a collection of **novel**, carefully designed **holistic** Arabic language benchmarking evaluation datasets that address these critical limitations. AraEval serves as a native Arabic benchmark, ensuring cultural, linguistic, and normative alignment with Arabic-speaking communities. Our contributions include:

1. AraEval includes **24,378 novel** samples across knowledge, reasoning, truthfulness, and instruction-following (Table 1).
2. AraEval facilitates detailed diagnostic assessments of model performance, enabling the

identification of specific strengths and weaknesses in reasoning, instruction-following, and knowledge retention. (Figures 1, 3, 4 and 7 and tables 7 to 10)

3. AraEval includes higher Arabic token coverage than ArabicMMLU and OpenAI’s Arabic-translated MMMLU (Figure 5 and table 16).
4. AraEval supports both log-probability-based and API-based evaluation schemes, facilitating seamless assessment of both open and close-source models.

## 2 AraEval Evaluation Suite

We contribute seven datasets of Arabic benchmarks, which vary in capabilities as shown in Table 1.

Task	Type	Dataset	Test Split	Dev Split
Knowledge	MCQ	AraPro	5001	110
Knowledge	MCQ	IEN MCQ	9990	190
Knowledge	Boolean	IEN TF	5823	190
Reasoning	MCQ	AraMath	605	5
Reasoning	MCQ	ETEC	1887	5
Instruction following	Generation	AraIFEval	536	-
Truthfulness	MCQ	AraTruthfulQA	536	5
Total			24,378	

Table 1: AraEval tasks splits statistics.

### 2.1 Design Principles

To establish a comprehensive Arabic benchmark for evaluating LLMs across diverse tasks, we developed our datasets based on the following principles:

**Human-curated or human-validated:** Every dataset of AraEval is meticulously created by experts or rigorously validated by humans to ensure the highest standards of quality and relevance. This guarantees that the questions, answers, and annotations are both accurate and meaningful, reflecting real-world scenarios and challenges. The validation criteria were task-specific, and human validators received specialized training on the respective tasks before beginning the validation process. The validation process was conducted by three humans where majority agreement was taken as the final verdict.

**Granularity for fine-grained evaluation:** Our datasets are designed with a high level of granularity, enabling detailed evaluation and nuanced insights into model performance. Fine-grained labels allow for the analysis of specific areas of strength and weakness, making the datasets particularly useful for diagnostic and comparative studies.

**Cultural and normative alignment:** All datasets are thoughtfully aligned with Arabic cul-

ture, values, and norms. This ensures the content is appropriate, contextually relevant, and reflective of the diverse realities of Arabic-speaking communities, allowing for more authentic and reliable evaluations.

## 2.2 Datasets Overview

### 2.2.1 AraPro

This dataset comprises 5,001 multiple-choice questions (MCQs) carefully crafted by university professors across 19 distinct knowledge domains. These experts were selected and instructed to create MCQs that reflect the competencies expected of professionals in their respective fields. Therefore, the questions evaluate LLMs in achieving professional-level competency within these domains. A detailed breakdown of the knowledge domains and the corresponding number of questions is provided in Table 10, while we show subject categories distribution in Figure 7.

### 2.2.2 IEN

The global pandemic of COVID-19 has challenged the world and inevitably the education sector. In Saudi Arabia, the Ministry of Education responded by launching the IEN<sup>2</sup> platform as part of its broader e-learning and distance education strategy.

The IEN platform includes a vast repository of more than 1.5 million questions and answers, meticulously classified into varying levels of difficulty. This extensive database not only supports differentiated learning, but also enables customized assessments that address the unique needs and abilities of students at every stage of their educational journey.

A representative subset that covers all grades, subjects and levels of difficulty was randomly selected from the IEN platform as shown in Table 1, the selection contains 5,823 samples as true/false questions and 9,990 MCQs. Figure 1 shows the detailed distributions of the questions and subjects per grade level. Table 8 and Table 9 provide more granular details about the dataset.

### 2.2.3 AraMath

AraMath consists of 605 MCQs derived from AraMath (Alghamdi et al., 2022), which includes mathematical word problems, and the solution is an equation that solves the problem. We reformulated the dataset and converted it to a multiple-choice problem (MCQ). The correct answer is extracted from the equation by parsing the formulas, and

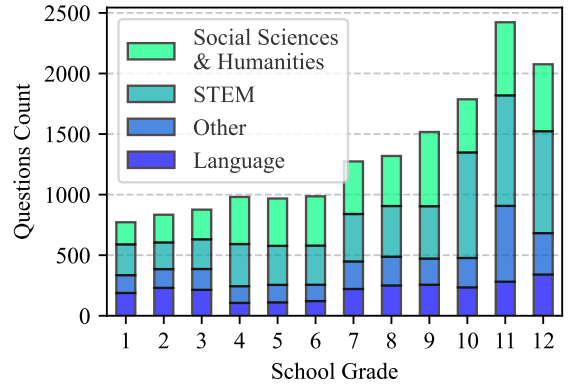


Figure 1: Course and grade level coverage for TF and MCQ IEN datasets combined.

three random distractors were generated to complete the set of options.

Human annotators meticulously reviewed and validated the dataset to ensure the accuracy of the equations in representing the mathematical word problems. They also assessed choice distinctiveness, verifying that all answer choices were unique and free of duplicates, and answer correctness, ensuring that the labeled answer corresponded to the correct choice.

### 2.2.4 ETEC

The Education & Training Evaluation Commission (ETEC)<sup>3</sup> serves as an independent regulatory body responsible for evaluating, measuring and accrediting qualifications in education and training in both the public and private sectors in Saudi Arabia. Its role includes ensuring and enhancing the quality and efficiency of educational and training institutions, programs, and their outcomes. The commission offers more than 42 types of qualification tests spanning all educational levels from K12 to professional levels. A subset of 1887 MCQs were chosen from different types of tests that include: a) Qudurati: A series of tests offered to students from 3rd grade elementary school to 10th grade to assess their level of general aptitude in comprehension, analysis, reasoning, and application, focusing on their readiness for learning. b) Professional Educational Occupation License Test: A standardized assessment tool to measure applicants' competency in general and specialized educational standards for on-the-job teachers.

<sup>2</sup><https://ien.edu.sa/>

<sup>3</sup><https://etec.gov.sa/en/>

### 2.2.5 AraIFEval

AraIFEval is an Arabic instruction-following (IF) evaluation benchmark designed to automatically assess language models’ compliance with specified instructions through verifiable methods. The dataset consists of 535 instances, each containing two to four verifiable instructions that can be validated using deterministic programming approaches. An example of the AraIFEval dataset with verifiable instructions is shown in Appendix D.3.

We created a collection of 23 Arabic verifiable instructions, inspired by Zhou et al. (2023). To construct the dataset, we randomly selected open questions from our data to serve as seed prompts. We generated IF prompts by randomly combining two to four instructions for each prompt, carefully ensuring logical consistency and avoiding contradictions between instructions. The dataset was then reviewed by humans for quality assurance. The Arabic verifiable instructions are presented in Appendix E, while the dataset distribution is detailed in Figure 8. To enable automatic response verification, we implemented regex-based category phrase checking. We followed Zhou et al. (Zhou et al., 2023) evaluation approach to assess instruction-following capabilities following strict and loose criteria. Similar to Fourier et al. (2024), we only report strict accuracy in this work.

### 2.2.6 AraTruthfulQA

Inspired by TruthfulQA Lin et al. (2021), this Arabic benchmark evaluates the truthfulness of LLM-generated responses to questions designed to elicit common misconceptions. The benchmark targets questions that some individuals may answer incorrectly due to false beliefs or misinformation. It comprises carefully curated questions spanning diverse categories, including religion, science, and nutrition, with a particular emphasis on prevalent misconceptions in the Arab world.

To ensure cultural and contextual relevance, we carefully reviewed the original TruthfulQA dataset and selected 287 questions that align with Arabic cultural norms and beliefs. These questions were translated into Arabic by human experts to preserve accuracy and nuance. Additionally, we crafted 249 culturally relevant questions of similar complexity and depth, specifically addressing common misconceptions in the Arab world, further enhancing the benchmark’s comprehensiveness.

## 3 Experiments

### 3.1 Setup

In this paper, we integrate the AraEval benchmark with the LM Evaluation Harness framework (Gao et al., 2024). We evaluate both open-source and closed models in zero-shot and few-shot settings, utilizing the test and dev sets; except for AraIFEval, where only zero-shot results are reported. To mitigate the token bias issue (Alzahrani et al., 2024), we have ensured a balanced distribution of the correct answer’s position in the MCQs datasets that have four choices such as AraMath, ETEC, and AraPro (see Figure 9). In the fewshots setting for IEN-MCQs, IEN-TF, and AraPro, we selected the few-shot examples that match the domain of the target question, in order to reduce the impact of out-of-domain questions in the few-shot samples.

#### 3.1.1 Open Models Setup

In order to evaluate the open-source models, since we can access their weights, we computed the log-probability for the choices in the MCQ datasets and reported the normalized accuracy. We used labels (A, B, C, D, etc.) to calculate log probabilities, except for AraTruthfulQA, where we calculated the log-probability of the choice label followed by the context of the choices.

As for AraIFEval, it was implemented as a generation task in LM Evaluation Harness where we report both prompt and instruction strict accuracies.

#### 3.1.2 Closed Models Setup

To evaluate the closed-source models for the AraEval suite, we implemented a generation-based evaluation using the LM Evaluation Harness framework. Since closed models can only be accessed through APIs and do not provide token-level probabilities (logprobs), we adapted all benchmark tasks in AraEval to a generation-based format to suit such models. We set the generation temperature to 0.0 to ensure consistency and determinism in the model responses. The closed-sourced models evaluated over this setting include GPT-4o (Hurst and et al., 2024), Claude-3.5-Sonnet (Anthropic, 2024), and Gemini-1.5-Pro (Gemini, 2024). For the multiple-choice tasks, such as the IEN datasets MCQ and TF, ETEC, AraMath, AraPro, and AraTruthfulQA, we applied filters that extract the model’s selected answer from its generated response. Such filters ensure that the extracted response corresponds exactly to one of the pro-



vided answer choices. After processing the model outputs, accuracy was calculated by comparing the extracted responses to the gold-standard labels using an exact match criterion.

### 3.2 Baselines

We evaluate a range of Arabic and state-of-the-art multilingual models to assess the utility of our evaluation suite. To this end, we design a series of experiments that: (1) compare model performance across various tasks, analyzing fine-grained results across different domains, (2) examine knowledge retention across different model sizes within the same family, and (3) compare base and instruct (chat) models to assess their relative strengths. Our evaluation covers models shown in Table 6, considering variants with 7B, 13B, 30B, and 70+B parameters to study scaling trends and performance variations.

### 3.3 Results

Zero-shot results for instruct models are shown in Table 2, while zero-shot for base models and five-shot results of base and instruct models are presented in Table 3, Table 4, and Table 5, respectively, in Appendix A. We report normalized accuracy for all tasks, and similar to [Fourrier et al. \(2024\)](#), we report strict prompt-level and instruction-level accuracy [Zhou et al. \(2023\)](#) for AraIFEval. Strict prompt-level accuracy refers to the percentage of prompts in which all verifiable instructions are correctly followed, while strict instruction-level accuracy is the percentage of verifiable instructions that are successfully followed.

The results reveal notable performance variations across models, model sizes, and shot settings. GPT4o, Claude, and Gemini demonstrate the highest performance across most tasks, consistently outperforming other models. Qwen 32B and 72B models and ALLaM 30B follow closely, showing robust performance across multiple tasks, especially in IEN MCQs and IEN TF. Llama 70B performs well but lags behind top-tier models, particularly in reasoning and advanced knowledge tasks including ETEC, AraPro, and AraMath, where its scores remain in the high 60s to low 70s. Among the Arabic models, these tasks remain challenging to Jais-family models where they underperform, while the AceGPT 32B model demonstrates improved performance; however, it falls short of achieving 70% accuracy.

The impact of model scaling varies across different types of tasks. For example, AraMath shows the most significant improvements with scaling, where Qwen 7B achieves an accuracy of 71.24% that increases to 92.07% with Qwen 32B. Similarly, Llama 3.3 70B achieves 69.92% compared to 32.73% with Llama 3.1 8B. Conversely, AraTruthfulQA do not exhibit the same level of improvement. For example, the Qwen models—7B, 14B, and 72B—achieve comparable accuracy rates of 52.8%, 58.4%, and 57.84%, respectively, while the Qwen 32B model outperforms them slightly with a higher accuracy of 61.19%.

The results highlight distinct patterns in task difficulty levels. Certain tasks, such as IEN MCQ and IEN TF, demonstrate consistently high accuracy across multiple models, suggesting a lower level of difficulty. This outcome is expected, as these tasks primarily consist of questions covering K01–K12 school subjects, which involve fundamental concepts and factual recall, making them easier for language models to handle. Other advanced knowledge and reasoning tasks, such as ETEC, AraPro, and AraMath, show a wider variance in scores, highlighting higher difficulty level. For ETEC, performance varies significantly across models, with Claude Sonnet 3.5 (85.9%) and Gemini Pro 1.5 (83.31%) achieving high scores, but Llama 8B is struggling at 45.68%. Similar trends are seen in AraMath and AraPro, where high variance is observed across models, with GPT4o achieving 81.16% and 80.86%, respectively, and Llama 8B scoring 32.73% and 52.51%, respectively. AraIFEval exhibit consistently low performance across all model families, indicating inherent difficulty. Even the strongest models achieve relatively low scores, compared to other tasks, with Claude sonnet 3.5 at 53.73%.

Most models benefit from few-shot prompting, but the degree of improvement varies. For instance, Qwen models show substantial improvements, particularly Qwen 7B, which gains over 10% in IEN MCQ, while Jais-family models struggle with few-shot prompting, with Jais-13B experiencing a performance drop in ETEC from 48.65% to 26.76%. Instruct models consistently outperform base models, particularly in AraMath, AraIFEval, and AraTruthfulQA. For example, Qwen 72B-Instruct scores 87.51% on AraIFEval, while its base counterpart achieves only 50.31%, highlighting the impact of instruction tuning on instruction following. Similarly, in AraTruthfulQA, ALLaM

34B Instruct scores 81.53%, whereas its base version achieves 64.18%, showing that fine-tuning improves truthfulness and misinformation resistance. However, for simpler knowledge-based tasks like IEN MCQ, the gap is smaller. In some cases, base models outperform their instruct counterparts, as seen in IEN MCQ, where Qwen 72B Base scores 90.77%, surpassing the 86.77% of its instruct version. Few-shot prompting benefits base models more than instruct models, as seen in the AraMath task, where Qwen 72B improves from 88.60% (0-shot) to 95.87% (5-shot). Overall, instruction tuning significantly enhances reasoning, alignment, and reliability, while larger base models still perform well in factual retrieval.

## 4 Analysis

### 4.1 Cross-Models Analysis

AraEval aggregates 7 datasets into a single score representing general Arabic capabilities. Inspired by [Fourrier et al. \(2024\)](#), we take the average normalized score across benchmarks, which is defined as:

$$\text{Norm. Score} = 100 \cdot \frac{\text{Raw Score} - \text{Baseline}}{100 - \text{Baseline}} \quad (1)$$

This transformation assigns a normalized score of 0% for the random baseline and 100% for a perfect score, with the rest linearly interpolated. In effect, this unifies score variances across benchmarks; It increases the contribution of benchmarks with high random baselines, such as true/false benchmarks, such that their scores span  $[0, 100]$  instead of  $[50, 100]$ . The final score is the mean of the 7 normalized benchmark scores. Five-shot evaluation is used whenever applicable to decouple formatting from base model evaluation.

Figure 2 illustrates the relationship between model size and AraEval accuracy for several prominent model families, including Qwen 2.5, Llama 3, Jais Family, AceGPT v2, ALLaM, and ALLaM Adapted. Across all model families, there is a consistent trend of increasing accuracy as model size scales from 7B to 70B parameters. This suggests that larger models are better equipped to capture the complexities of the Arabic language, benefiting from richer parameterization. While all models demonstrated performance gains with increased size, ALLaM Base exhibited the most significant improvements, particularly in the small-to-mid size range (7B–30B), indicating the effectiveness of its architecture and training data for

Arabic-specific tasks. The sensitivity of AraEval to variations in model scale—from 7B to 70B parameters—further highlights the benchmark’s robustness. It effectively captures nuanced performance differences, making it particularly well-suited for fine-grained comparisons across diverse model configurations.

Although performance generally improved with size, diminishing returns became apparent beyond the 30B parameter mark for Qwen2.5 and for ALLaM instruct scaling from 7B to 30B. For these models, the accuracy gains were marginal compared to the more substantial improvements observed when scaling from 7B to 30B in Llama 3 instruct and ALLaM base. This suggests potential saturation points where further parameter increases yield limited benefits. This ability to detect performance plateaus is critical for guiding model scaling decisions and optimizing resource allocation.

Instruct models consistently outperform their Base counterparts across all size categories, underscoring the benchmark’s ability to reflect improvements from fine-tuning strategies aimed at aligning models with user instructions.

### 4.2 Fine-Grained Analysis

While average evaluation metrics provide a general overview of LLMs performance, fine-grained assessments offer deeper insights into specific capabilities and areas needing improvement. This detailed evaluation is crucial for understanding the strengths and weaknesses of LLMs in various contexts. Several approaches were proposed to reveal the fine-grained capabilities of models. FAC<sup>2</sup>E ([Wang et al., 2024a](#)) proposed a framework for better understanding LLM capabilities by dissociating Language and Cognitive capabilities allowing for a more detailed analysis of LLM performance. Similarly, the "FLASK" ([Ye et al., 2024](#)) evaluation protocol decomposes overall scoring into specific skill sets for each instruction, providing a fine-grained evaluation that enhances interpretability and reliability. To this extent, AraEval benchmark offers a deeper insight into the capabilities of LLMs by pinpointing model scoring not only at an overall view but more deeper such as grade, subject, and difficulty level, See Figure 3, 4 and 6. The variations in the figures indicate that the models performances varies and provide insightful remarks about how each model performs when compared to others, and at the same time will identify the gap or the deficiencies the model might

Model	IEN		AraPro	AraMath	ETEC	AraTruthfulQA	AraIFEval	
	MCQ	TF					Prompt	Instruction
ALLaM 7B-Instruct	93.10	83.14	73.51	70.08	70.38	71.83	59.51	82.46
Llama-3.1-8B-Instruct	59.23	71.73	52.51	32.73	45.68	54.29	53.36	79.32
Qwen2.5-7B-Instruct	66.38	78.46	64.63	71.24	64.12	52.8	28.17	65.19
ALLaM Adapted 13B-Instruct	<b>93.44</b>	83.75	74.69	78.68	73.87	67.16	59.33	83.14
Jais-family-13B-chat	62.95	68.68	57.53	42.64	48.65	56.53	17.16	54.27
Qwen2.5-14B-Instruct	80.51	77.64	69.11	80.17	72.18	58.4	68.66	86.76
ALLaM 34B-Instruct	93.29	86.83	79.52	60.50	74.24	78.36	67.16	86.76
AceGPT-v2-32B-chat	81.60	80.35	67.19	64.13	64.81	65.11	25.75	63.41
Jais-family-30B-16k-chat	74.88	68.76	62.79	50.74	53.31	63.99	16.60	54.95
Jais-family-30B-8k-chat	72.76	70.65	61.27	42.64	53.52	62.69	16.79	54.68
Qwen2.5-32B-Instruct	84.93	81.92	71.81	92.07	78.33	61.19	56.90	82.87
ALLaM Adapted 70B-Instruct	92.56	85.56	75.82	73.22	76.21	81.72	65.49	85.39
Jais-adapted-70B-chat	74.51	76.47	64.59	50.74	56.81	71.46	27.05	65.05
Llama-3.3-70B-Instruct	79.60	78.81	70.49	69.92	68.84	67.16	70.90	88.60
Qwen2.5-72B-Instruct	86.88	86.62	74.69	89.26	78.70	57.84	67.72	87.51
GPT-4o	92.03	88.97	80.86	81.16	79.39	87.69	70.90	88.12
Gemini pro 1.5	88.28	85.44	76.22	<b>96.36</b>	83.31	88.43	<b>74.81</b>	<b>90.17</b>
Claude Sonnet 3.5	86.17	<b>89.42</b>	<b>81.46</b>	88.6	<b>85.9</b>	<b>90.67</b>	53.73	80.14
Random baseline	30.77	50	25	25	25	23.46	0	0

Table 2: Overall results of instruct models across all AraEval benchmarks 0-shot.

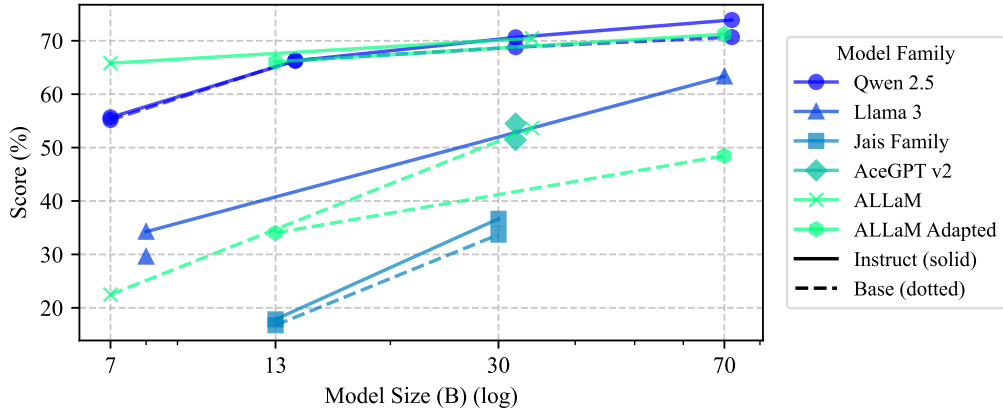


Figure 2: LLMs performance on AraEval for various model sizes. Instruct models are in solid lines, while Base models are in dashed lines.

suffer from. In Figure 3(b) it is noticeable that in the questions “Above average” there is more variance between the models compared to other types - “Average” or “Below average” difficulty questions. Similarly subjects like “Language” Figure 3(c), and “Humanities” (Figure 4) show similar trends where the performance of the models varies widely. Such nuances and observations are useful and insightful and reflect the utility of a high quality benchmark.

### 4.3 Vocabulary Coverage Analysis

A robust evaluation of large language models in Arabic requires not only challenging tasks, but also a comprehensive vocabulary coverage. In this work, we assess the vocabulary coverage of several

models across the Arabic datasets within our proposed benchmark AraEval, and compare it against Arabic MMLU (Koto et al., 2024) and OpenAI MMMLU (translated to Arabic) (OpenAI, 2024) two widely used benchmarks in the community.

As shown in Figure 5, the vocabulary coverage values are averaged across all models. AraEval achieves 74.05% coverage of Arabic tokens, closely aligning with OpenAI Arabic MMMLU (74.17%), while surpassing Arabic MMLU (66.38%). This coverage ensures that AraEval incorporates a diverse range of Arabic tokens, including domain-specific tokens from science, history, and literature.

This rich token representation makes AraEval a

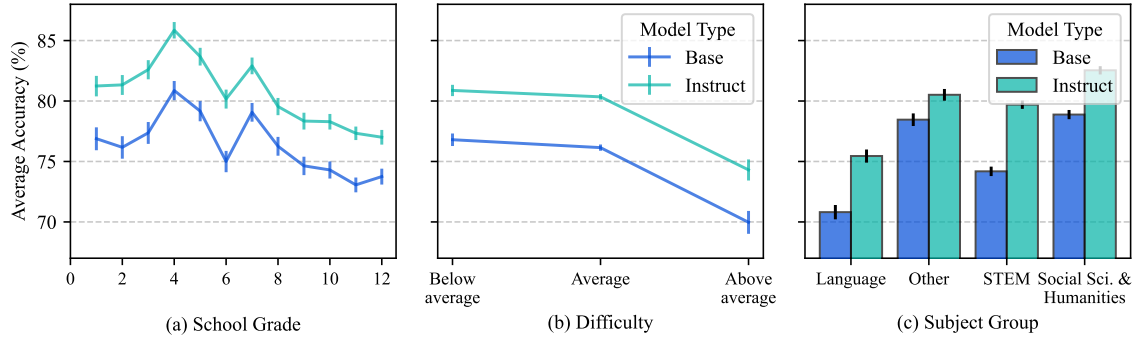


Figure 3: Average accuracies on all evaluated models for various IEN MCQ subsets. Error bars represent 95% confidence intervals of the average accuracy across all models.

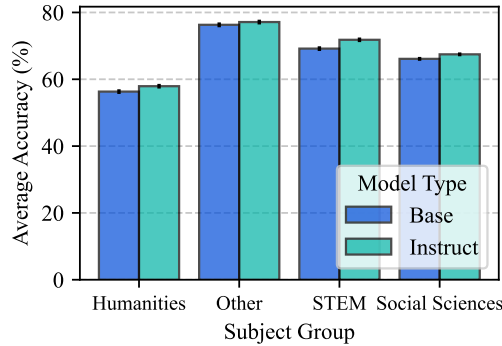


Figure 4: Average accuracies on all evaluated models for various AraPro subsets. Error bars represent 95% confidence intervals for the average across models.

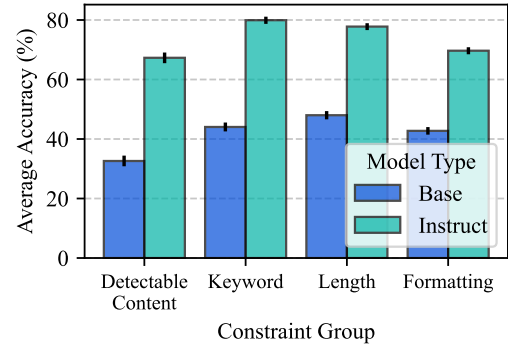


Figure 6: Average accuracies on all evaluated models for various AraFEval constraint subsets. Error bars represent 95% confidence intervals across models.

more faithful and challenging benchmark for evaluating LLM performance in Arabic. A detailed breakdown of the vocabulary coverage is provided in Table 16.

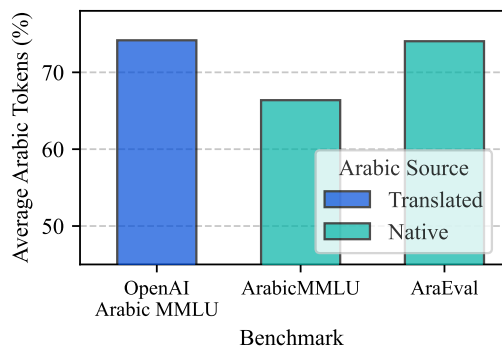


Figure 5: Average Arabic vocabulary coverage across various tokenizers. Details are presented in Table 16. AraEval covers a large portion of Arabic vocabulary without using translated data.

## 5 Conclusion

In this paper, we introduced AraEval, a comprehensive benchmark designed to rigorously evaluate the

advanced knowledge, reasoning, and instruction-following capabilities of foundation models within the Arabic context. Our evaluation highlights the robustness and diversity of the datasets within AraEval, offering key insights into their effectiveness in distinguishing model capabilities. Tasks like AraMath, AraPro, ETEC, and AraFEval prove highly challenging, effectively differentiating models based on their reasoning and problem-solving skills, making them strong indicators of true model competency. AraTruthfulQA effectively measures a model’s susceptibility to misinformation, revealing clear differences in truthfulness across models. Conversely, datasets such as IEN MCQ and IEN TF capture less advanced knowledge that some base models can handle. These findings emphasize the value of AraEval as a benchmarking tool for Arabic LLMs. The diversity of tasks ensures that models are tested across multiple dimensions of knowledge and reasoning, allowing for precise performance diagnostics. As the field of LLMs continues to evolve, AraEval provides a strong foundation for future evaluations, paving the way for more targeted advancements in Arabic NLP.



## 6 Limitations

Despite AraEval’s contribution to addressing the gap in comprehensive assessment datasets, several limitations warrant consideration. First, the dataset’s reliance on multiple-choice questions (MCQ) and true/false formats inherently constrains the evaluation of language models’ capabilities. These structured response formats may not adequately assess deeper levels of comprehension or the ability to generate creative solutions that more closely align with real-world applications. Second, the dataset’s predominant focus on Saudi curriculum introduces potential cultural bias. This geographical and cultural specificity may limit the dataset’s generalizability to educational contexts in other regions and cultures, potentially overlooking important cultural nuances and educational approaches from diverse educational systems. Third, the current benchmark’s scope is limited to text-based assessments, excluding evaluation capabilities for multi-modal models. This limitation becomes particularly significant as artificial intelligence increasingly requires the ability to process and synthesize information across various modalities, including visual, auditory, and textual data. These limitations suggest opportunities for future work to develop more comprehensive evaluation frameworks that incorporate open-ended responses, diverse cultural perspectives, and multi-modal assessment capabilities.

## 7 Ethical Considerations

All authors of this work acknowledge and adhere to the ACL Code of Ethics, upholding its principles throughout the research process. All domain experts and annotators involved in the creation and review of the datasets are official employees, who are fairly compensated based on mutually agreed-upon wage standards and working hours. These employment agreements fully comply with local labor regulations. Furthermore, we prioritize clear communication about how data and annotations are utilized, obtaining informed consent from domain experts and annotators before incorporating their contributions into our research. We are also dedicated to safeguarding their privacy throughout the annotation and data creation process, fostering an ethical and respectful research environment.

## References

- Ummar Abbas, Mohammad Shahmeer Ahmad, Firoj Alam, Enes Altinisik, Ehsannedin Asgari, Yazan Boshmaf, Sabri Boughorbel, Sanjay Chawla, Shammur Chowdhury, et al. 2025. Fanar: An arabic-centric multimodal generative ai platform. *arXiv preprint arXiv:2501.13944*.
- Saja AL-Tawalbeh and Mohammad Al-Smadi. 2020. [A benchmark arabic dataset for commonsense explanation](#). *ArXiv*, abs/2012.10251.
- Reem Alghamdi, Zhenwen Liang, and Xiangliang Zhang. 2022. [Armath: a dataset for solving arabic math word problems](#). In *Proceedings of the Language Resources and Evaluation Conference*, pages 351–362, Marseille, France. European Language Resources Association.
- Ebtesam Almazrouei, Ruxandra Cojocaru, Michele Baldo, Quentin Malartic, Hamza Alobeidli, Daniele Mazzotta, Guilherme Penedo, Giulia Campesan, Murgariya Farooq, Maitha Alhammedi, Julien Launay, and Badreddine Noune. 2023. [AlGhafa evaluation benchmark for Arabic language models](#). In *Proceedings of ArabicNLP 2023*, pages 244–275, Singapore (Hybrid). Association for Computational Linguistics.
- AlphaCode Team. 2023. [Alphacode 2 technical report](#). *Google DeepMind*.
- Norah Alzahrani, Hisham Alyahya, Yazeed Alnumay, Sultan AlRashed, Shaykhah Alsubaie, Yousef Almushayqih, Faisal Mirza, Nouf Alotaibi, Nora Al-Twairish, Areeb Alowisheq, M Saiful Bari, and Haidar Khan. 2024. [When benchmarks are targets: Revealing the sensitivity of large language model leaderboards](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13787–13805, Bangkok, Thailand. Association for Computational Linguistics.
- Anthropic. 2022. [The claude 3 model family: Opus, sonnet, haiku](#).
- AI Anthropic. 2024. Claude 3.5 sonnet model card addendum. *Claude-3.5 Model Card*, 3:6.
- Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. 2023. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. *arXiv preprint arXiv:2308.16884*.
- M Saiful Bari, Yazeed Alnumay, Norah A. Alzahrani, Nouf M. Alotaibi, Hisham A. Alyahya, Sultan Al-Rashed, Faisal A. Mirza, Shaykhah Z. Alsubaie, Hassan A. Alahmed, Ghadah Alabduljabbar, Raghad Alkhathran, Yousef Almushayqih, Raneem Alnajim, Salman Alsubaihi, Maryam Al Mansour, Majed Al-rubaian, Ali Alammari, Zaki Alawami, Abdulmohsen Al-Thubaity, Ahmed Abdelali, Jeril Kuriakose, Abdalghani Abujabal, Nora Al-Twairish, Areeb Alowisheq, and Haidar Khan. 2024. [Allam: Large language models for arabic and english](#). *Preprint*, arXiv:2407.15390.

- Yuri Chervonyi, Trieu H Trinh, Miroslav Olšák, Xiaomeng Yang, Hoang Nguyen, Marcelo Menegali, Junehyuk Jung, Vikas Verma, Quoc V Le, and Thang Luong. 2025. Gold-medalist performance in solving olympiad geometry with alphageometry2. *arXiv preprint arXiv:2502.03544*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2020. *Ethnologue: Languages of the World*, 23 edition. SIL International, Dallas.
- Ali El Filali, Hamza Alobeidli, Clémentine Fourrier, Basma El Amel Boussaha, Ruxandra Cojocaru, Nathan Habib, and Hakim Hacid. 2024. Open arabic llm leaderboard. <https://huggingface.co/spaces/OALL/Open-Arabic-LLM-Leaderboard>.
- Aaron Grattafiori et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. 2024. Open llm leaderboard v2. [https://huggingface.co/spaces/open-llm-leaderboard/open\\_llm\\_leaderboard](https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard).
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. [A framework for few-shot language model evaluation](#).
- Team Gemini. 2024. [Gemini 1.5: Unlocking multi-modal understanding across millions of tokens of context](#). *Preprint*, arXiv:2403.05530.
- Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caroline Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, et al. 2024. Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai. *arXiv preprint arXiv:2411.04872*.
- Paul Gómez-Rodríguez, Carlos Williams. 2023. [A confederacy of models: a comprehensive evaluation of LLMs on creative writing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14504–14528, Singapore. Association for Computational Linguistics.
- Google. 2024. [Gemini: A family of highly capable multimodal models](#). *Preprint*, arXiv:2312.11805.
- Google. 2025. Google gemini: Deep research. <https://blog.google/products/gemini/>.
- [google-gemini-deep-research/](https://google-gemini-deep-research/). Accessed: 2025-02-11.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Kunal Handa, Alex Tamkin, Miles McCain, Saffron Huang, Esin Durmus, Sarah Heck, Jared Mueller, Jerry Hong, Stuart Ritchie, Tim Belonax, Kevin K. Troy, Dario Amodei, Jared Kaplan, Jack Clark, and Deep Ganguli. 2025. [Which economic tasks are performed with ai? evidence from millions of claude conversations](#). *Anthropic*.
- Momchil Hardalov, Todor Mihaylov, Dimitrina Zlatkova, Yoan Dinkov, Ivan Koychev, and Preslav Nakov. 2020. EXAMS: A multi-subject high school examinations dataset for cross-lingual and multilingual question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5427–5444, Online. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Huang Huang, Fei Yu, Jianqing Zhu, Xuening Sun, Hao Cheng, Dingjie Song, Zhihong Chen, Abdulmohsen Alharthi, Bang An, Juncai He, et al. 2023. Acegpt, localizing large language models in arabic. *arXiv preprint arXiv:2309.12053*.
- Aaron Hurst and et al. 2024. [Gpt-4o system card](#). *Preprint*, arXiv:2410.21276.
- Inception. 2024. [Jais family model card](#). *Hugging Face*.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. [The state and fate of linguistic diversity and inclusion in the NLP world](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.
- Daniel Kahneman. 2011. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, New York.
- Mohammad Abdullah Matin Khan, M Saiful Bari, Do Long, Weishi Wang, Md Rizwan Parvez, and Shafiq Joty. 2024. XCodeEval: An execution-based large scale multilingual multitask benchmark for code understanding, generation, translation and retrieval. In *Proceedings of the 62nd Annual Meeting*

829	<i>of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6766–6805, Bangkok, Thailand. Association for Computational Linguistics.	890
830		891
831		892
832	Fajri Koto, Haonan Li, Sara Shatnawi, Jad Doughman, Abdelrahman Sadallah, Aisha Alraeesi, Khalid Al-mubarak, Zaid Alyafeai, Neha Sengupta, Shady Shehata, Nizar Habash, Preslav Nakov, and Timothy Baldwin. 2024. <a href="#">ArabicMMLU: Assessing massive multitask language understanding in Arabic</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 5622–5640, Bangkok, Thailand. Association for Computational Linguistics.	893
833		894
834		895
835		896
836		897
837		898
838		899
839		900
840		901
841	Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan Rossi, and Thien Nguyen. 2023. <a href="#">Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 318–327, Singapore. Association for Computational Linguistics.	902
842		903
843		904
844		905
845		906
846		907
847		908
848		909
849		910
850	LAION-AI. 2025. Open-assistant: A chat-based assistant for task understanding and dynamic interaction. <a href="https://github.com/LAION-AI/Open-Assistant">https://github.com/LAION-AI/Open-Assistant</a> . Accessed: 2025-02-12.	911
851		912
852		913
853		914
854	Md Tahmid Rahman Laskar, Sawsan Alqahtani, M Saiful Bari, Mizanur Rahman, Mohammad Abdullah Matin Khan, Haidar Khan, Israt Jahan, Amran Bhuiyan, Chee Wei Tan, Md Rizwan Parvez, Enamul Hoque, Shafiq Joty, and Jimmy Huang. 2024. <a href="#">A systematic survey and critical review on evaluating large language models: Challenges, limitations, and recommendations</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 13785–13816, Miami, Florida, USA. Association for Computational Linguistics.	915
855		916
856		917
857		918
858		919
859		920
860		921
861		922
862		923
863		924
864		925
865	Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Huang. 2023. <a href="#">A systematic study and comprehensive evaluation of ChatGPT on benchmark datasets</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 431–469, Toronto, Canada. Association for Computational Linguistics.	926
866		927
867		928
868		929
869		930
870		931
871		932
872		933
873	Latent Space. 2024. In the arena: How lmsys changed llm benchmarking forever. <a href="https://www.latent.space/p/lmarena">https://www.latent.space/p/lmarena</a> . Blog post (accessed: 11 February 2025).	934
874		935
875		936
876		937
877	Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models. In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22</i> , Red Hook, NY, USA. Curran Associates Inc.	938
878		939
879		940
880		941
881		942
882		943
883		944
884		945
885		946
886	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2024. <a href="#">CMMLU: Measuring massive multitask language understanding in Chinese</a> . In <i>Findings of</i>	947
887		948
888		949
889		950
		951
	<i>the Association for Computational Linguistics: ACL 2024</i> , pages 11260–11285, Bangkok, Thailand. Association for Computational Linguistics.	
	Juhao Liang, Zhenyang Cai, Jianqing Zhu, Huang Huang, Kewei Zong, Bang An, Mosen Alharthi, Juncai He, Lian Zhang, Haizhou Li, Benyou Wang, and Jinchao Xu. 2024. <a href="#">Alignment at pre-training! towards native alignment for arabic LLMs</a> . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. <i>arXiv preprint arXiv:2109.07958</i> .	
	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024. Deepseek-v3 technical report. <i>arXiv preprint arXiv:2412.19437</i> .	
	Holy Lovenia, Rahmad Mahendra, Salsabil Maulana Akbar, Lester James V. Miranda, Jennifer Santoso, Elyanah Aco, Akhdan Fadhlil, Jonibek Mansurov, Joseph Marvin Imperial, Onno P. Kampman, Joel Ruben Antony Moniz, Muhammad Ravi Shulthan Habibi, Frederikus Hudi, Railey Montalan, Ryan Ignatius, Joanito Agili Lopo, William Nixon, Börje F. Karlsson, James Jaya, Ryandito Diandaru, Yuze Gao, Patrick Amadeus, Bin Wang, Jan Christian Blaise Cruz, Chenxi Whitehouse, Ivan Halim Parmonangan, Maria Khelli, Wenyu Zhang, Lucky Susanto, Reynard Adha Ryanda, Sonny Lazuardi Hermawan, Dan John Velasco, Muhammad Dehan Al Kautsar, Willy Fitra Hendria, Yasmin Moslem, Noah Flynn, Muhammad Farid Adilazuarda, Haochen Li, Johanes Lee, R. Damanhuri, Shuo Sun, Muhammad Reza Qorib, Amirbek Djanibekov, Wei Qi Leong, Quyet V. Do, Niklas Muennighoff, Tanrada Pansuwan, Ilham Firdausi Putra, Yan Xu, Ngee Chia Tai, Ayu Purwarianti, Sebastian Ruder, William Tjhi, Peerat Limkonchotiwat, Alham Fikri Aji, Sedrick Keh, Genta Indra Winata, Ruochen Zhang, Fajri Koto, Zheng-Xin Yong, and Samuel Cahyawijaya. 2024. <a href="#">Seacrowd: A multilingual multimodal data hub and benchmark suite for southeast asian languages</a> . <i>arXiv preprint arXiv: 2406.10118</i> .	
	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz	



952	Ferrandis, Lingming Zhang, Sean Hughes, Thomas	David Rein, Betty Li Hou, Asa Cooper Stickland, Jack-	1010
953	Wolf, Arjun Guha, Leandro von Werra, and Harm	son Petty, Richard Yuanzhe Pang, Julien Dirani, Ju-	1011
954	de Vries. 2024. <a href="#">Starcoder 2 and the stack v2: The</a>	lian Michael, and Samuel R Bowman. 2023. Gpqa: A	1012
955	<a href="#">next generation</a> . <i>Preprint</i> , arXiv:2402.19173.	graduate-level google-proof q&a benchmark. <i>arXiv</i>	1013
956	Grégoire Mialon, Clémentine Fourrier, Craig Swift,	<i>preprint arXiv:2311.12022</i> .	1014
957	Thomas Wolf, Yann LeCun, and Thomas Scialom.	Yusuke Sakai, Hidetaka Kamigaito, and Taro Watanabe.	1015
958	2023. Gaia: a benchmark for general ai assistants.	2024. <a href="#">mCSQA: Multilingual commonsense reason-</a>	1016
959	<i>arXiv preprint arXiv:2311.12983</i> .	<a href="#">ing dataset with unified creation strategy by language</a>	1017
960	Mistral. 2024. <a href="#">Au large</a> .	<a href="#">models and humans</a> . In <i>Findings of the Association</i>	1018
961	Basel Mousi, Nadir Durrani, Fatema Ahmad, Md Arif	<i>for Computational Linguistics: ACL 2024</i> , pages	1019
962	Hasan, Maram Hasanain, Tameem Kabbani, Fahim	14182–14214, Bangkok, Thailand. Association for	1020
963	Dalvi, Shammur Absar Chowdhury, and Firoj	Computational Linguistics.	1021
964	Alam. 2024. Aradice: Benchmarks for dialectal	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,	1022
965	and cultural capabilities in llms. <i>arXiv preprint</i>	Satheesh Katipomu, Haonan Li, Fajri Koto, William	1023
966	<i>arXiv:2409.11404</i> .	Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming	1024
967	Ahmad Mustapha, Hadi Al-Khansa, Hadi Al-Mubasher,	Chen, Osama Mohammed Afzal, Samta Kamboj,	1025
968	Aya Mourad, Ranam Hamoud, Hasan El-Husseini,	Onkar Pandit, Rahul Pal, Lalit Pradhan, Zain Muham-	1026
969	Marwah Al-Sakkaf, and Mariette Awad. 2024.	mad Mujahid, Massa Baali, Xudong Han, Son-	1027
970	<a href="#">Arastem: A native arabic multiple choice question</a>	dos Mahmoud Bsharat, Alham Fikri Aji, Zhiqiang	1028
971	<a href="#">benchmark for evaluating llms knowledge in stem</a>	Shen, Zhengzhong Liu, Natalia Vassilieva, Joel Hes-	1029
972	<a href="#">subjects</a> . <i>Preprint</i> , arXiv:2501.00559.	tness, Andy Hock, Andrew Feldman, Jonathan Lee,	1030
973	OpenAI. 2022. <a href="#">Chatgpt: Optimizing language models</a>	Andrew Jackson, Hector Xuguang Ren, Preslav	1031
974	<a href="#">for dialogue</a> .	Nakov, Timothy Baldwin, and Eric Xing. 2023a.	1032
975	OpenAI. 2024. <a href="#">Multilingual massive multitask lan-</a>	<a href="#">Jais and jais-chat: Arabic-centric foundation and</a>	1033
976	<a href="#">guage understanding (mmmlu)</a> .	<a href="#">instruction-tuned open generative large language</a>	1034
977	OpenAI. 2025. <a href="#">Introducing deep re-</a>	<a href="#">models</a> . <i>Preprint</i> , arXiv:2308.16149.	1035
978	<a href="#">search</a> . <a href="https://openai.com/index/introducing-deep-research/">https://openai.com/index/</a>	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,	1036
979	<a href="#">introducing-deep-research/</a> . Accessed:	Satheesh Katipomu, Haonan Li, Fajri Koto, William	1037
980	2025-02-11.	Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming	1038
981	OpenAI. 2025. simple-evals. <a href="https://github.com/openai/simple-evals">https://github.com/</a>	Chen, et al. 2023b. <a href="#">Jais and jais-chat: Arabic-</a>	1039
982	<a href="#">openai/simple-evals</a> . Accessed: 2025-02-11.	<a href="#">centric foundation and instruction-tuned open ge-</a>	1040
983	Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li,	<a href="#">nerative large language models</a> . <i>arXiv preprint</i>	1041
984	Josephina Hu, Hugh Zhang, Sean Shi, Michael Choi,	<i>arXiv:2308.16149</i> .	1042
985	Anish Agrawal, Arnab Chopra, et al. 2025. Human-	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang,	1043
986	ity’s last exam. <i>arXiv preprint arXiv:2501.14249</i> .	Suraj Srivats, Soroush Vosoughi, Hyung Won Chung,	1044
987	Edoardo Maria Ponti, Goran Glavaš, Olga Majewska,	Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das,	1045
988	Qianchu Liu, Ivan Vulić, and Anna Korhonen. 2020.	and Jason Wei. 2023. <a href="#">Language models are multi-</a>	1046
989	<a href="#">XCOPA: A multilingual dataset for causal common-</a>	<a href="#">lingual chain-of-thought reasoners</a> . In <i>The Eleventh</i>	1047
990	<a href="#">sense reasoning</a> . In <i>Proceedings of the 2020 Con-</i>	<i>International Conference on Learning Representa-</i>	1048
991	<i>ference on Empirical Methods in Natural Language</i>	<i>tions, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> .	1049
992	<i>Processing (EMNLP)</i> , pages 2362–2376, Online. As-	OpenReview.net.	1050
993	sociation for Computational Linguistics.	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mah-	1051
994	Zhaozhi Qian, Farooq Altam, Muhammad Alqurishi, and	davi, Jason Wei, Hyung Won Chung, Nathan Scales,	1052
995	Riad Souissi. 2024. <a href="#">Camelevel: Advancing cultur-</a>	Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl,	1053
996	<a href="#">ally aligned arabic language models and benchmarks</a> .	et al. 2023. Large language models encode clinical	1054
997	<i>Preprint</i> , arXiv:2409.12623.	knowledge. <i>Nature</i> , 620(7972):172–180.	1055
998	Qwen, :, An Yang, Baosong Yang, Beichen Zhang,	Guijin Son, Hanwool Lee, Sungdong Kim, Seungone	1056
999	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,	Kim, Niklas Muennighoff, Taekyoon Choi, Cheon-	1057
1000	Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin,	bok Park, Kang Min Yoo, and Stella Biderman. 2024.	1058
1001	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	<a href="#">Kmmlu: Measuring massive multitask language un-</a>	1059
1002	Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	<a href="#">derstanding in korean</a> . <i>Preprint</i> , arXiv:2402.11548.	1060
1003	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li,	Aarohi Srivastava and et al. 2023. <a href="#">Beyond the imitation</a>	1061
1004	Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji	<a href="#">game: Quantifying and extrapolating the capabilities</a>	1062
1005	Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang	<a href="#">of language models</a> . <i>Trans. Mach. Learn. Res.</i> , 2023.	1063
1006	Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	1064
1007	Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru	Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	1065
1008	Zhang, and Zihan Qiu. 2025. <a href="#">Qwen2.5 technical</a>	Adam R Brown, Adam Santoro, Aditya Gupta,	1066
1009	<a href="#">report</a> .	Adrià Garriga-Alonso, et al. 2022. <a href="#">Beyond the</a>	1067
		<a href="#">imitation game: Quantifying and extrapolating the</a>	1068
		<a href="#">capabilities of language models</a> . <i>arXiv preprint</i>	1069
		<i>arXiv:2206.04615</i> .	1070



1071	Qwen Team. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	1131
1072		1132
1073	Reka Team, Aitor Ormazabal, Che Zheng, Cyprien de Masson d’Autume, Dani Yogatama, Deyu Fu, Donovan Ong, Eric Chen, Eugenie Lamprecht, Hai Pham, et al. 2024. Reka core, flash, and edge: A series of powerful multimodal language models. <i>arXiv preprint arXiv:2404.12387</i> .	1133
1074		1134
1075		1135
1076		1136
1077		1137
1078		1138
1079	Xiaoqiang Wang, Lingfei Wu, Tengfei Ma, and Bang Liu. 2024a. <a href="#">FAC<sup>2</sup>E: Better understanding large language model capabilities by dissociating language and cognition</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 13228–13243, Miami, Florida, USA. Association for Computational Linguistics.	1139
1080		1140
1081		1141
1082		1142
1083		1143
1084		1144
1085		1145
1086	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. 2024b. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. <i>arXiv preprint arXiv:2406.01574</i> .	1146
1087		1147
1088		1148
1089		1149
1090		1150
1091		1151
1092	Zeyu Wang. 2024. <a href="#">CausalBench: A comprehensive benchmark for evaluating causal reasoning capabilities of large language models</a> . In <i>Proceedings of the 10th SIGHAN Workshop on Chinese Language Processing (SIGHAN-10)</i> , pages 143–151, Bangkok, Thailand. Association for Computational Linguistics.	1152
1093		1153
1094		
1095		
1096		
1097		
1098	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> .	
1099		
1100		
1101		
1102		
1103	BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Lucic, François Yvon, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. <i>arXiv preprint arXiv:2211.05100</i> .	
1104		
1105		
1106		
1107		
1108		
1109	Vikas Yadav, Steven Bethard, and Mihai Surdeanu. 2019. <a href="#">Quick and (not so) dirty: Unsupervised selection of justification sentences for multi-hop question answering</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2578–2589, Hong Kong, China. Association for Computational Linguistics.	
1110		
1111		
1112		
1113		
1114		
1115		
1116		
1117		
1118	Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. 2024. <a href="#">Flask: Fine-grained language model evaluation based on alignment skill sets</a> . <i>Preprint</i> , arXiv:2307.10928.	
1119		
1120		
1121		
1122		
1123	Zheng-Xin Yong, Cristina Menghini, and Stephen H Bach. 2023. Low-resource languages jailbreak gpt-4. <i>arXiv preprint arXiv:2310.02446</i> .	
1124		
1125		
1126	Arda Yüksel, Abdullatif Köksal, Lütfi Kerem Senel, Anna Korhonen, and Hinrich Schuetze. 2024. <a href="#">TurkishMMLU: Measuring massive multitask language understanding in Turkish</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 7035–7055, Miami, Florida, USA. Association for Computational Linguistics.	
1127		
1128		
1129		
1130		
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. <a href="#">HellaSwag: Can a machine really finish your sentence?</a> In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800, Florence, Italy. Association for Computational Linguistics.	
	Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. Available: <a href="https://github.com/google-research/google-research/tree/master/instruction_following_eval">https://github.com/google-research/google-research/tree/master/instruction_following_eval</a> .	
	Jianqing Zhu, Huang Huang, Zhihang Lin, Juhao Liang, Zhengyang Tang, Khalid Almubarak, Mosen Alharthi, Bang An, Juncai He, Xiangbo Wu, Fei Yu, Junying Chen, Zhuoheng Ma, Yuhao Du, Yan Hu, He Zhang, Emad A. Alghamdi, Lian Zhang, Ruoyu Sun, Haizhou Li, Benyou Wang, and Jinchao Xu. 2024. Second language (arabic) acquisition of llms via progressive vocabulary expansion. <i>arXiv preprint arXiv:2412.12310</i> .	

## A Additional Results

In addition to Table 2 in the main paper. We also show the remaining 0 and 5 shot base and instruct model results in Tables 3 to 5.

## B Evaluated Models

Table 6 outlines the LLMs used in our evaluation with additional details.

Size	Model	Creator	Access
7B	Qwen 2.5 (Qwen et al., 2025)	Alibaba	weights
8B	Llama 3.1 (et al., 2024)	Meta	weights
7B	ALLaM (Bari et al., 2024)	SDAIA	weights
14B	Qwen 2.5	Alibaba	weights
13B	Jais family 13b chat (Sengupta et al., 2023a; Inception, 2024)	InceptionAI	weights
13B	ALLaM Adapted	SDAIA	weights
32B	Qwen 2.5	Alibaba	weights
30B	Jais family 30b 8k-chat	InceptionAI	weights
30B	Jais family 30b 16k-chat	InceptionAI	weights
32B	AceGPT (Zhu et al., 2024; Liang et al., 2024)	FreedomIntelligence	weights
34B	ALLaM	SDAIA	weights
72B	Qwen 2.5	Alibaba	weights
70B	Llama 3.3	Meta	weights
70B	Jais-adapted 70b-chat	InceptionAI	weights
70B	ALLaM Adapted	SDAIA	weights
—	GPT4o (Hurst and et al., 2024)	OpenAI	API
—	Gemini pro 1.5 (Gemini, 2024)	Google	API
—	Claude 3.5 Sonnet (Anthropic, 2024)	Anthropic	API

Table 6: Instruct models considered

## C Related Work

Evaluating LLMs requires comprehensive benchmark datasets that assess knowledge, reasoning, and language understanding. These datasets can be categorized into general-purpose and domain-specific types, ensuring models are both broadly competent and specialized.

### C.1 General-Purpose Datasets

General-purpose datasets evaluate a model’s versatility across tasks like question-answering, translation, and commonsense reasoning. The Massive Multitask Language Understanding (MMLU) dataset (Hendrycks et al., 2021) measures general knowledge across 57 subjects, with adaptations for languages such as Korean (KMMLU) (Son et al., 2024), Turkish (TurkishMMLU) (Yüksel et al., 2024), and Chinese (CMMLU) (Li et al., 2024). OpenAI has also translated MMLU into 14 languages, including Arabic (OpenAI, 2024).

HellaSwag (Zellers et al., 2019) evaluates commonsense reasoning through multiple-choice questions, with multilingual extensions like XCOPA (Ponti et al., 2020) and mCSQA (Sakai et al., 2024). Grade School Math 8K (GSM8K) (Cobbe et al., 2021) focuses on quantitative reasoning, extended to ten languages via MGSM (Shi et al.,

2023). Finally, BigBench (Srivastava and et al., 2023) offers over 200 diverse tasks to test LLM capabilities across various domains.

### C.2 Domain-Specific Datasets

Domain-specific datasets evaluate LLMs in specialized fields. ARC-Challenge (Yadav et al., 2019) tests science reasoning, with Arabic versions like Okapi ARC-Challenge (Lai et al., 2023) and Al-Ghafa Evaluation Benchmark (Almazrouei et al., 2023). Minerva Math (Lewkowycz et al., 2022) assesses mathematical reasoning, while CausalBench (Wang, 2024) evaluates causal inference across textual, mathematical, and coding domains. Multi-MedQA (Singhal et al., 2023) combines six medical datasets to evaluate clinical knowledge, making it essential for healthcare-related tasks.

### C.3 Arabic Datasets

Few datasets have been explicitly developed to evaluate LLMs in Arabic, but recent efforts have made significant progress. One notable example is ArabicMMLU (Koto et al., 2024), a comprehensive multiple-choice question benchmark designed to assess reasoning and knowledge capabilities of LLMs in Modern Standard Arabic. Developed with input from native speakers across North Africa, the Levant, and the Gulf, it includes 14,575 questions spanning 40 diverse tasks. These tasks cover subjects such as STEM, social sciences, humanities, and the Arabic language, sourced from educational materials in various Arabic-speaking countries. The dataset reflects a range of educational levels.

Another important contribution is AraSTEM (Mustapha et al., 2024), which focuses on STEM subjects like mathematics, physics, chemistry, biology, computer science, and medicine. This dataset comprises multiple-choice questions sourced from elementary, secondary, and higher education levels, ensuring broad coverage of difficulty and topics. It was carefully compiled from multiple internet sources to ensure diversity and comprehensiveness.

Efforts to adapt existing English evaluation datasets for Arabic include the AlGhafa Arabic LLM Benchmark (Almazrouei et al., 2023). This benchmark consists of 11 datasets translated or modified from English benchmarks, verified by native Arabic speakers. Similarly, the Benchmark Arabic Dataset for Commonsense Explanation (AL-Tawalbeh and Al-Smadi, 2020) translates the original English ComVE task into Arabic. It contains 12,000 instances, each presenting an Ara-

Model	IEN		AraPro	AraMath	ETEC	AraTruthfulQA	AraIFEval	
	MCQ	TF					Prompt	Instruction
ALLaM 7B Base	58.83	57.53	49.41	20.33	39.43	44.78	3.73	29.56
Llama-3.1-8B	64.30	53.37	51.07	26.61	42.77	54.29	7.28	41.50
Qwen2.5-7B	77.10	77.21	61.75	67.93	59.62	71.08	6.72	44.57
ALLaM Adapted 13B Base	63.41	66.82	54.85	23.14	40.65	50	6.53	38.50
Jais-family-13B	38.04	53.61	31.15	31.90	28.40	50	6.90	40.75
Qwen2.5-14B	83.63	69.17	68.45	79.17	69.69	66.98	10.82	47.78
ALLaM 34B Base	83.49	57.05	72.71	48.10	62.43	53.54	<b>17.16</b>	<b>55.15</b>
AceGPT-v2-32B	78.49	65.81	65.85	54.71	58.77	63.81	8.02	45.26
Jais-family-30B-16k	67.03	54.42	54.29	28.10	42.13	48.88	11.01	45.12
Jais-family-30B-8k	58.76	60.90	55.21	26.12	42.82	48.13	11.57	48.74
Qwen2.5-32B	85.03	<b>82.05</b>	71.43	81.82	75.57	73.13	11.75	46.35
ALLaM Adapted 70B Base	75.76	75.49	64.19	35.54	54.90	59.33	3.17	24.30
Jais-adapted-70B	70.35	60.23	61.79	37.69	44.89	61.19	9.89	43.21
Qwen2.5-72B	<b>88.79</b>	79.75	<b>73.89</b>	<b>88.60</b>	<b>78.01</b>	<b>78.73</b>	14.93	50.31
Random baseline	30.77	50	25	25	25	23.46	0	0

Table 3: Overall results of base models across all AraEval benchmarks 0-shot.

Model	IEN		AraPro	AraMath	ETEC	AraTruthfulQA
	MCQ	TF				
ALLaM 7B Base	63.78	64.62	55.77	18.02	43.46	43.28
Llama-3.1-8B	71.22	62.56	59.29	39.67	47.96	51.49
Qwen2.5-7B	81.66	78.88	66.55	75.70	65.34	75.75
ALLaM Adapted 13B Base	72.62	71.29	62.93	23.47	50.98	59.70
Jais-family-13B	32.43	58.78	40.35	26.45	33.39	42.35
Qwen2.5-14B	86.54	83.77	72.53	92.56	75.68	83.96
ALLaM 34B Base	86.22	81.68	77.16	51.74	65.77	64.18
AceGPT-v2-32B	83.02	80.37	70.11	66.45	65.02	72.95
Jais-family-30B-16k	72.93	69.72	65.09	35.87	51.40	53.36
Jais-family-30B-8k	71.57	68.28	63.05	32.23	51.03	52.24
Qwen2.5-32B	87.95	<b>86.02</b>	74.99	94.05	79.65	82.28
ALLaM Adapted 70B Base	83.04	76.83	72.45	48.26	63.01	79.48
Jais-adapted-70B	78.33	74.36	66.97	51.24	52.20	77.24
Qwen2.5-72B	<b>90.77</b>	85.35	<b>77.86</b>	<b>95.87</b>	<b>82.25</b>	<b>84.33</b>
Random baseline	30.77	50	25	25	25	23.46

Table 4: Overall results of base models across all AraEval benchmarks 5-shot

bic sentence that defies commonsense, accompanied by three explanatory options. The task is to identify the best explanation for why the sentence is nonsensical.

Qian et al. (2024) introduced CamelEval, a suite of three test sets designed to evaluate general instruction following, factuality, and cultural alignment in Arabic. Each test set includes 805 carefully curated cases reflecting the nuances of the Arabic language and culture.

While these datasets significantly advance the evaluation of Arabic LLMs, they also exhibit certain limitations. For instance, ArabicMMLU and AraSTEM may not fully capture the diversity of educational systems, cultural nuances, and histori-

cal contexts across Arabic-speaking countries. Despite sourcing questions from multiple regions, ArabicMMLU might struggle to encompass the full spectrum of curricula and perspectives in the Arab world. Similarly, AraSTEM, while focusing on STEM subjects, may not adequately represent the varied educational strategies and cultural contexts found in different Arabic-speaking nations.

Additionally, translating English datasets into Arabic, such as in the case of AlGhafa and the Benchmark Arabic Dataset for Commonsense Explanation, presents challenges. Translations may fail to preserve cultural nuances and contextual meanings inherent in the original language, leading to potential misinterpretations. Furthermore, these

Model	IEN		AraPro	AraMath	ETEC	AraTruthfulQA
	MCQ	TF				
ALLaM 7B	92.61	84.36	73.97	73.06	70.06	71.46
Llama-3.1-8B-Instruct	65.15	59.92	57.45	35.70	47.75	58.58
Qwen2.5-7B-Instruct	78.18	77.98	65.97	71.74	64.92	69.96
ALLaM Adapted 13B	92.51	83.03	74.93	75.04	73.40	70.34
Jais-family-13B-chat	53.65	60.24	32.99	26.61	26.76	48.69
Qwen2.5-14B-Instruct	81.10	80.37	71.31	82.81	73.29	70.34
ALLaM 34B	<b>93.00</b>	87.65	80.70	62.81	73.87	81.53
AceGPT-v2-32B-chat	82.98	73.28	68.23	64.46	65.77	67.54
Jais-family-30B-16k-chat	71.43	64.14	62.57	41.49	49.28	61.75
Jais-family-30B-8k-chat	67.40	71.90	60.61	33.39	45.52	59.7
Qwen2.5-32B-Instruct	85.17	82.83	73.45	91.90	78.01	76.12
ALLaM Adapted 70B	92.22	85.08	76.74	74.88	76.10	84.14
Jais-adapted-70B-chat	77.34	76.64	68.23	45.62	57.50	77.43
Llama-3.3-70B-Instruct	81.27	80.01	72.53	70.91	67.89	70.71
Qwen2.5-72B-Instruct	86.77	86.74	75.66	92.89	79.12	71.27
GPT-4o	91.70	89.64	81.46	83.47	79.49	90.11
Gemini pro 1.5	84.06	87.09	78.28	<b>94.88</b>	84.31	84.14
Claude Sonnet 3.5	88.6	<b>90.74</b>	<b>83.96</b>	79.83	<b>86.43</b>	<b>93.47</b>
Random baseline	30.77	50	25	25	25	23.46

Table 5: Overall results of instruct models across all AraEval benchmarks 5-shot.

datasets may not align well with the educational curricula and cultural contexts of Arabic-speaking countries, where educational systems and cultural norms vary significantly. This misalignment can result in evaluations that do not accurately reflect the capabilities of Arabic-centric LLMs in real-world applications.

## D AraEval Datasets

In this section, we detail each dataset used in AraEval, including fine-grained analyses, task statistics, and example samples.

### D.1 Domain and Subject Distribution

Table 8 and Table 9 show distribution for both IEN MCQ and IEN TF, respectively, in terms of study stage, difficulty level, and subjects.

AraPro subjects distribution is presented in Table 10 and category distribution in Figure 7. For AraIFEval, we show the distribution of constraint groups in Figure 8, while Table 7 shows the distribution of instructions, where each sample comprises multiple instructions.

### D.2 MCQ Datasets Distribution

Figure 9 shows the options distribution in AraEval datasets.

Category	Count	Percent (%)
number words at least	265	18.09
number paragraphs	225	15.36
response language	139	9.49
title	135	9.22
keyword frequency	135	9.22
number words at most	87	5.94
include keywords	63	4.30
forbidden words	60	4.10
number bullets	48	3.28
letter frequency	46	3.14
postscript	34	2.32
first word in i-th paragraph	33	2.25
check end	27	1.84
number sentences at least	25	1.71
minimum number highlighted section	22	1.50
json format	21	1.43
multiple sections	20	1.37
quotation	20	1.37
number placeholder	14	0.96
repeat prompt	13	0.89
two responses	12	0.82
number sentences at most	11	0.75
no commas	10	0.68
<b>Total</b>	<b>1465</b>	<b>-</b>

Table 7: Category distribution and percentage of AraIFEval dataset.

### D.3 Dataset Examples

Figure 10 illustrates the construction of verifiable instructions in AraIFEval: the upper part shows the original (normal) instruction, while the bottom part shows the instruction after adding verifiable prompts.



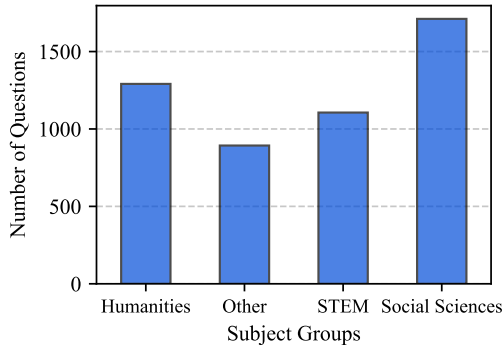


Figure 7: Subject distribution of AraPro.

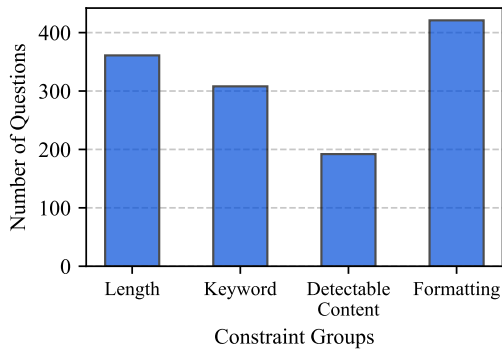


Figure 8: Constraint distribution of AraIFEval.

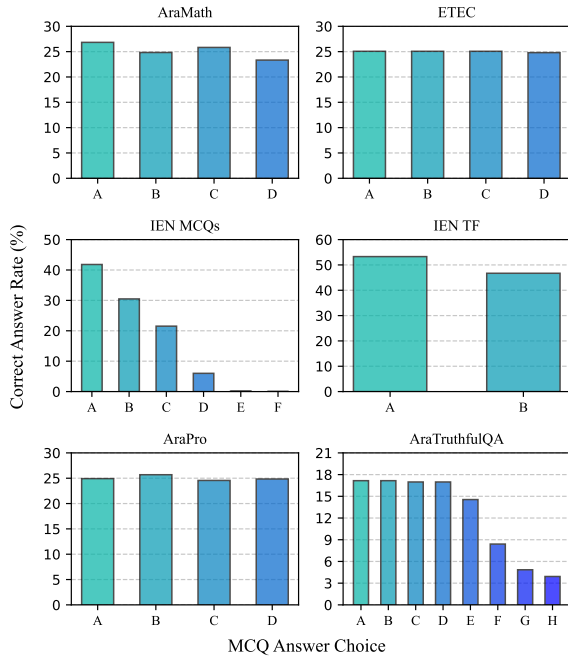


Figure 9: Distribution percentage of the correct answer in each MCQ dataset of AraEval.

## E AraIFEval Prompts

Table 11 shows the instructions categories prompts in AraIFEval.

Category	#Subject/Specialty	#Questions
<i>In terms of study stages</i>		
Secondary education	17	3747
Primary education	10	3739
Intermediate education	11	2504
<i>In terms of difficulty level</i>		
Easy	17	1834
Medium	17	7505
Hard	17	651
<i>In terms of Levels</i>		
K01	8	551
K02	8	583
K03	8	595
K04	9	680
K05	9	660
K06	9	670
K07	10	769
K08	10	892
K09	11	906
K10	13	1057
K11	13	1293
K12	13	1240
<i>Breakdown by Subject/Specialty</i>		
Social Studies and National Ed	—	844
Biology	—	178
Research and Information Sour	—	92
Family and Health Education	—	854
Physical Education	—	517
Art Education	—	829
Computer Science	—	1003
Mathematics	—	799
Science	—	944
Administrative Sciences	—	284
Islamic Studies	—	1209
Behavioral Sciences	—	267
Physics	—	239
Chemistry	—	220
English Language	—	637
Arabic Language	—	980
Environmental Science	—	93
<b>Total</b>	<b>17</b>	<b>9990</b>

Table 8: Statistics of IEN MCQs.

Prompt

من هي الشركة المسؤولة عن استخراج الذهب في المملكة العربية السعودية؟

**Buckwalter:** mn hy Al\$rkp Alms&wlp En AstxrAj Al\*hb fy Almmklp AlErbypp AlsEwdyp ?

**English Translation:** Who is the company responsible for extracting gold in the Kingdom of Saudi Arabia?

Prompt with Verifiable Instruction

من هي الشركة المسؤولة عن استخراج الذهب في المملكة العربية السعودية؟ يجب أن يحتوي ردك على ما لا يقل عن 50 كلمة. يجب أن تظهر الكلمة "الذهب" في ردك ثلاث مرات على الأقل. يجب أن يكون ردك بالكامل باللغة العربية.

**Buckwalter:** mn hy Al\$rkp Alms&wlp En AstxrAj Al\*hb fy Almmklp AlErbypp AlsEwdyp? yjb >n yHtwy rdk EIY mAlA yql En 50 klmp. yjb >n tZhr klmp "Al\*hb" fy rdk vAv mrAt EIY Al>ql. yjb >n ykwn rdk bAlkAml bAllgp AlErbypp.

**English Translation:** Who is the company responsible for extracting gold in the Kingdom of Saudi Arabia? Your response must contain at least 50 words. The word 'gold' must appear in your response at least three times. Your response must be entirely in Arabic.

Figure 10: Example of verifiable instruction created of an existing instruction in Arabic.

Category	#Subject/Specialty	#Questions
<i>In terms of study stages</i>		
Secondary education	17	2539
Primary education	10	1678
Intermediate education	11	1606
<i>In terms of difficulty level</i>		
Easy	17	1360
Medium	17	4195
Hard	17	268
<i>In terms of levels</i>		
K01	8	221
K02	8	251
K03	8	281
K04	9	301
K05	9	308
K06	9	316
K07	10	505
K08	11	490
K09	11	611
K10	13	730
K11	13	973
K12	13	836
<i>Breakdown by Subject/Specialty</i>		
Social Studies and Nation	–	482
Biology	–	159
Research and Information	–	99
Family and Health Educat.	–	453
Physical Education	–	421
Art Education	–	380
Computer Science	–	598
Mathematics	–	507
Science	–	421
Administrative Sciences	–	161
Islamic Studies	–	558
Behavioral Sciences	–	233
Physics	–	133
Chemistry	–	197
English Language	–	394
Arabic Language	–	530
Environmental Science	–	97
<b>Total</b>	<b>17</b>	<b>5823</b>

Table 9: Statistics of IEN TF.

## F Dataset Curation and Validation

The guidelines for domain experts on creating AraPro can be found in Table 12, while the validation guidelines for AraMath are presented in Table 13. The guideline for validation of AraIFEval is detailed in Table 14, and the guidelines for AraTruthfulQA are provided in Table 15.

## G Tokenizer Vocabulary Coverage

Table 16 shows the models’ vocabulary coverage across the Arabic datasets within AraEval compared to MMLU and OpenAI MMLU benchmarks.

Subject	#Question
<i>Breakdown by Subject/Specialty</i>	
Sociology	403
Biology	212
Management	197
Arabic Literature	558
Economics	397
History	297
Computing	199
Religion	299
Sports	396
Mathematics	200
Politics	414
Physics	97
Chemistry	200
Arabic Linguistics	434
Finance	100
Human Resources	200
Engineering	98
Psychology	200
Earth Sciences	100
<b>Total</b>	<b>5001</b>

Table 10: Statistics of AraPro.

## H GPU Time

GPU time for running evaluation on AraEval datasets is reported in Table 17.

Dataset	7B	13B	30B	70B
AraPro (0 shot)	447.65	969.77	4326.20	9770.33
AraPro (5 shot)	328.78	576.82	1434.85	2459.53
IEN MCQ (0 shot)	420.02	463.81	1268.43	2129.42
IEN MCQ (5 shot)	552.10	867.71	2875.39	4196.97
IEN TF (0 shot)	269.64	357.27	1232.53	1686.52
IEN TF (5 shot)	321.30	514.43	1344.34	2677.28
AraMath (0 shot)	44.55	62.17	1676.55	3623.28
AraMath (5 shot)	61.19	94.08	253.83	396.62
ETEC (0 shot)	153.76	172.00	351.70	550.40
ETEC (5 shot)	226.07	367.63	1031.75	1685.91
AraIFEval (0 shot)	7051.31	6954.25	29382.06	29724.12
AraTruthfulQA (0 shot)	514.21	844.75	4443.01	9924.95
AraTruthfulQA (5 shot)	250.30	494.59	1226.33	2111.18

Table 17: GPU time for different model sizes. The reported time is in seconds and is the average across all models of the corresponding size.

Instruction Category	Prompt
include_keywords	قم بتضمين الكلمات المفتاحية في ردك qm btDmyn AlklmAt AlmfAtAHyp (keyword1), (keyword2) fy rdk. Include the keywords (keyword1) and (keyword2) in your response
keyword_frequency	يجب أن تظهر الكلمة (word) في ردك (N) مرة yjb >n tZhr Alklmp (word) fy rdk (N) mrp The word (word) must appear in your response (N) times
forbidden_words	لا تتم بتضمين الكلمات المحظورة IA tqm btDmyn AlklmAt AlmHZwrp Do not include the (forbidden word)
letter_frequency	يجب أن يظهر الحرف (letter) في ردك (N) مرة yjb >n yZhr AlHrf (letter ) fy rdk (N) mrp. The letter (letter) must appear (N) times in your response
response_language	يجب أن يكون ردك بالكامل باللغة (language) ولا يُسمح بأي لغة أخرى yjb >n ykwn rdk bAlkAml bAllgp (language) wIA ysmH blgp >xrY Your response must be entirely in (language), and no other language is allowed
number_paragraphs	يجب أن يحتوي ردك على عدد معين من الفقرات yjb >n yHtwy rdk ElY Edd mEyn mn AlfqrAt Your response must contain (N) paragraphs
number_words_at_least	أجب بما لا يقل عن (N) كلمة >jb bmA IA yql En (N) klmp Answer with at least (N) words
number_words_at_most	أجب بما لا يزيد عن (N) كلمة >jb bmA lAyzyd En (N) klmp Answer with (N) words at most
number_sentences_at_least	أجب بما لا يقل عن (N) جملة >jb bmA IA yql En (N) jmlp Answer with at least (N) sentences
number_sentences_at_most	أجب بما لا يزيد عن (N) جملة >jb bmA IA yzyd En (N) jmlp Answer with (N) sentences at most
first_word_in_i-th_paragraph	يجب أن نحوي الإجابة على عدد معين من الفقرات وتبدأ إحدى الفقرات بكلمة محددة yjb >n tHwy Al<jAbp ElY Edd mEyn mn AlfqrAt wtbd> <HdY AlfqrAt bkmlp mHddp ) The answer must contain a specific number of paragraphs, with one of the paragraphs starting with a specific word
postscript	يرجى إضافة ملاحظة توضيحية في نهاية ردك تبدأ ب (postscript marker) yrjy <DAfp mlAHZp twDyHyp fy nhAyp rdk tbd> b (postscript marker) Please add a clarifying note at the end of your response, starting with (postscript marker)
number_placeholder	يجب أن يحوي ردك على عدد من مواضع الترميز تمثل بأقواس مربعة yjb >n yHwy rdk Ely Edd mn mwADE Altmyz tmvl b>qwAs mrbEp Your response must contain at least (N) placeholders, represented using square brackets
number_bullets	يجب أن يحتوي ردك على عدد معين من النقاط yjb >n yHtwy rdk ElY Edd mEyn mn AlnqAT Your response must contain a specific number of points.
title	يجب أن يحتوي ردك على عنوان بين أقواس مزدوجة yjb >n yHtwy rdk ElY EnwAn byn >qwAs mzdwpj Your response must include a title enclosed in double angle brackets
minimum_number_high-lighted_section	قم بتسليط الضوء على عدد م من الأقسام على الأقل qm btslyT AIdw' ElY Edd m mn Al>qsAm ElY Al>qI Highlight at least Highlight at least sections.
multiple_sections	يجب أن يحتوي ردك على عدد م من الأقسام. ضع علامة على بداية كل قسم yjb >n yHtwy rdk ElY Edd m mn Al>qsAm . DE ElAmp ElY bdAyp kl qsm Your response must contain N sections. Place a section separator at the beginning of each section
json_format	يجب أن يكون الرد بالكامل بتنسيق JSON yjb >n ykwn Alrd bAlkAml btnsyq JSON Your response must be entirely formatted in JSON
repeat_prompt	كرر المدخل دون تغيير ثم قدم إجابتك krr Almdxl dwn tgyyr vm qdm <jAbtk Repeat the input without modification then respond to the prompt
two_responses	قدم إجابتين مختلفتين. الردود فقط يجب فصلها ب ٦ رموز نجوم qdm <jAbtn mxtlftyn. Alrdwd fqT yjb fSlhA b 6 rmwz njwm Provide two different answers. The responses should only be separated by six asterisk symbols
end_checker	انه ردك بالعبارة المحددة Anh rdk bAlEbArp AlmHddp End your response with specific phrase
quotation	يجب أن يكون ردك بالكامل بين علامات اقتباس مزدوجة yjb >n ykwn rdk byn ElAmAt AqtbAs mzdwpj Your response should be between double quotation mark
no-comma	امتنع عن استخدام فواصل في ردك AmtnE En AstxdAm fwASl fy rdk Don't use comma in your response

Table 11: Instructions categories prompts. We used buckwalter transliteration to transliterate Arabic instructions.

Section	Guidelines
<b>Objective</b>	The goal of these MCQs is to evaluate Large Language Models (LLMs) in achieving professional-level competency in your field of expertise. Each question should reflect real-world knowledge, critical thinking, and problem-solving skills relevant to industry standards. The data you create will only be used for research purposes.
<b>Question Structure</b>	Each MCQ should consist of: <ul style="list-style-type: none"> <li>• A clear and concise question that assesses knowledge, application, or analysis.</li> <li>• Four answer choices (A, B, C, D), with only one correct answer.</li> </ul>
<b>Guidelines for Crafting Questions</b>	<ul style="list-style-type: none"> <li>• Ensure relevance to key competencies in the profession.</li> <li>• Avoid ambiguity, excessive complexity, or unnecessary jargon.</li> <li>• Use practical scenarios, case studies, or problem-solving situations where possible.</li> <li>• Maintain a mix of basic, intermediate, and advanced questions.</li> <li>• Avoid testing trivial facts; focus on meaningful concepts.</li> </ul>
<b>Answer Choices</b>	<ul style="list-style-type: none"> <li>• One clear correct answer that is indisputably accurate.</li> <li>• Three plausible distractors that are incorrect but not obviously wrong.</li> </ul>
<b>Example Question Format</b>	<p><b>Question:</b> What is the primary purpose of risk assessment in cybersecurity?</p> <ul style="list-style-type: none"> <li>• A) To eliminate all potential threats</li> <li>• B) To identify, analyze, and mitigate security risks</li> <li>• C) To ensure compliance with industry regulations only</li> <li>• D) To monitor network traffic for suspicious activity</li> </ul> <p><b>Correct Answer:</b> B) To identify, analyze, and mitigate security risks  <b>Domain:</b> Computing  <b>Difficulty level:</b> Intermediate</p>
<b>Submission Format</b>	<ul style="list-style-type: none"> <li>• Provide questions in a structured format (Question, Options, Correct Answer, Domain, Difficulty level).</li> <li>• Ensure accuracy and relevance.</li> <li>• Submit questions in a spreadsheet as instructed.</li> </ul>
<b>Review Process</b>	All questions will be reviewed for accuracy, clarity, and alignment with professional competencies before finalization.

Table 12: Guidelines for Creating AraPro Dataset.



Section	Guidelines
<b>Objective</b>	The purpose of this validation process is to ensure the accuracy, consistency, and quality of a dataset containing mathematical word problems. Annotators are responsible for verifying the correctness of equations, answer choices, and labels to maintain data integrity. This dataset is used to evaluate mathematical reasoning capability of Large Language Models (LLMs). The data will be used for research purposes only.
<b>Dataset Components</b>	Each data entry consists of: <ul style="list-style-type: none"> <li>- <b>Mathematical Word Problem:</b> A problem statement requiring mathematical reasoning.</li> <li>- <b>Equation:</b> The corresponding mathematical equation representing the problem.</li> <li>- <b>Answer Choices (A, B, C, D):</b> Four distinct answer options.</li> <li>- <b>Correct Answer:</b> The solution to the problem.</li> <li>- <b>Answer Label:</b> The letter (A, B, C, or D) corresponding to the correct choice.</li> </ul>
<b>Validation Criteria</b>	<p><b>1. Accuracy of Equations</b></p> <ul style="list-style-type: none"> <li>- Verify that the equation correctly represents the given word problem.</li> <li>- Ensure the mathematical formulation aligns with the intended logic.</li> <li>- Check for errors in mathematical symbols, operations, and missing components.</li> </ul> <p><b>2. Choice Distinctiveness</b></p> <ul style="list-style-type: none"> <li>- Confirm that all four answer choices are unique and do not repeat.</li> <li>- Ensure that distractor options are plausible but incorrect.</li> <li>- Avoid choices that are too similar (e.g., minor rounding differences).</li> </ul> <p><b>3. Answer Correctness</b></p> <ul style="list-style-type: none"> <li>- Solve the problem independently and compare it with the provided correct answer.</li> <li>- Cross-check that the correct answer matches the labeled answer choice.</li> <li>- If errors are found, provide corrected answers and labels.</li> </ul> <p><b>4. Presence of Correct Answer</b></p> <ul style="list-style-type: none"> <li>- Ensure that the correct answer is one of the four given choices.</li> <li>- If the correct answer is missing from the options, flag the entry for correction.</li> </ul> <p><b>5. Formatting and Consistency</b></p> <ul style="list-style-type: none"> <li>- Ensure uniform formatting across all dataset entries.</li> <li>- Verify that symbols, units, and mathematical notation follow standard conventions.</li> </ul> <p><b>6. Logical Soundness</b></p> <ul style="list-style-type: none"> <li>- Assess whether the problem makes sense mathematically and linguistically.</li> <li>- Check for unintended biases or misleading wording.</li> </ul>
<b>Annotation Process</b>	<ol style="list-style-type: none"> <li>1. Read the problem statement carefully and understand its context.</li> <li>2. Examine the provided equation and ensure it correctly models the problem.</li> <li>3. Verify that the correct answer is calculated accurately.</li> <li>4. Confirm that all answer choices are unique and logically reasonable.</li> <li>5. Check that the correct answer exists within the four given choices.</li> <li>6. Cross-check the labeled answer against the correct answer.</li> <li>7. If discrepancies are found, document corrections and flag the entry for review.</li> </ol>
<b>Error Reporting &amp; Corrections</b>	Annotators should log any errors found, specifying: <ul style="list-style-type: none"> <li>- <b>Entry ID:</b> The unique identifier of the dataset entry.</li> <li>- <b>Issue Type:</b> (Equation Error, Answer Mismatch, Duplicate Choices, Missing Correct Answer, Formatting Issue, etc.).</li> <li>- <b>Correction:</b> The revised equation, answer choice, or label.</li> <li>- <b>Comments:</b> Additional notes explaining the error.</li> </ul>
<b>Final Review &amp; Approval</b>	<ul style="list-style-type: none"> <li>- After validation, a second-level review may be conducted to ensure error-free dataset entries.</li> <li>- Approved entries will be included in the final dataset, while flagged entries undergo correction and re-evaluation.</li> </ul>

Table 13: Guidelines for Human Annotators to validate AraMath Dataset.

Section	Guidelines
<b>Objective</b>	The purpose of this task is to ensure that each instance in this data accurately represents its instructed prompt and instruction categories. Annotators review the dataset for logical consistency, completeness, and correctness. This dataset is used to evaluate instruction following capability of Large Language Models (LLMs). The data will be used for research purposes only.
<b>Dataset Components</b>	Each data entry consists of: - <b>Instructed Prompt:</b> A textual prompt containing verifiable instructions. - <b>Instruction Categories:</b> A set of verifiable instructions used in the prompt.
<b>Validation Criteria</b>	<b>1. Contradiction Check</b> - Ensure that no contradictory instructions exist within the instructed prompt. - Flag instances where conflicting instructions lead to logical inconsistencies.  <b>2. Instruction Completeness</b> - Verify that all instruction categories in the instruction set are explicitly mentioned in the instructed prompt. - If an instruction is missing, annotate it as an omission.  <b>3. Prompt Coverage</b> - Ensure that all instructions present in the instructed prompt are correctly identified in the instruction set. - If additional, unlisted instructions are found, flag them for review.  <b>4. Logical Coherence</b> - Assess whether the prompt flows naturally and follows a coherent structure. - Check for redundant, unclear, or ambiguous wording.  <b>5. Formatting and Standardization</b> - Verify that instruction labels and categories follow the predefined taxonomy. - Ensure proper punctuation, spelling, and grammar for clarity.
<b>Annotation Process</b>	1. Read the instructed prompt carefully to understand its structure and intent. 2. Compare the instruction categories with the prompt to check for completeness. 3. Identify and flag any contradictory instructions within the prompt. 4. Verify that no instruction is missing from the instruction set. 5. Ensure that no extra, unlisted instructions are present in the prompt. 6. Check for formatting, clarity, and coherence issues. 7. Document errors and suggest corrections where necessary.
<b>Error Reporting &amp; Corrections</b>	Annotators should log errors using the following details: - <b>Entry ID:</b> Unique identifier of the dataset instance. - <b>Issue Type:</b> (Contradiction, Missing Instruction, Extra Instruction, Formatting Issue, etc.). - <b>Correction:</b> Suggested revision for the prompt or instruction set. - <b>Comments:</b> Additional explanation of the issue.
<b>Final Review &amp; Approval</b>	- A second-level review may be conducted to ensure high-quality validation. - Approved entries are included in the final dataset, while flagged entries undergo correction and re-evaluation.

Table 14: Guidelines for Validation of AraIFEval Dataset.

Section	Guidelines
<b>Objective</b>	The purpose of this review process is to evaluate multiple-choice questions (MCQs) for alignment with Arabic cultural norms and beliefs, ensuring that the content is appropriate, respectful, and contextually relevant. Additionally, validated MCQs should be translated into Arabic while maintaining their original meaning and intent. This dataset is used to evaluate truthfulness of Large Language Models (LLMs). The data will be used for research purposes only.
<b>Dataset Components</b>	Each MCQ consists of: - <b>Question:</b> The main stem of the MCQ. - <b>Four Answer Choices:</b> Options (A, B, C, D), with only one correct answer. - <b>Correct Answer Label:</b> The letter corresponding to the correct answer.
<b>Validation Criteria</b>	<b>1. Cultural Alignment</b> - Ensure that the question and answer choices do not conflict with Arabic cultural and social values. - Avoid topics that may be considered sensitive or inappropriate in an Arabic cultural context. - Verify that examples, names, and scenarios used in the MCQ are relevant and culturally recognizable.  <b>4. Translation Guidelines</b> - Translate only the MCQs that align with Arabic cultural norms. - Maintain the original intent and meaning of the question while using culturally appropriate phrasing. - Adapt idiomatic expressions or region-specific references to ensure clarity for Arabic speakers. - Use Modern Standard Arabic (MSA) for translation, avoiding dialect-specific terms.
<b>Annotation Process</b>	1. Read the MCQ carefully and assess its cultural appropriateness. 2. If the MCQ is <b>not aligned</b> , flag it and provide a justification. 3. If the MCQ is <b>aligned</b> , proceed with translation while preserving accuracy and clarity. 4. Ensure that all answer choices remain meaningful and distinguishable after translation. 5. Verify that the correct answer remains unchanged in meaning. 6. Document any modifications made during translation for transparency.

Table 15: Guidelines for Reviewing and Translating TruthfulQA dataset.

Benchmark	ALLaM-7B	ALLaM-34B	ALLaM-Adapted	Jais-Family	Jais-Adapted	Qwen-2.5*	Llama-3**
AraFEval	7.80	7.54	9.72	6.64	8.98	37.29	35.79
ETEC	32.37	33.34	38.10	28.39	35.53	67.22	58.74
IEN MCQs	53.64	56.15	60.22	48.33	56.82	77.34	63.36
IEN TF	36.24	36.84	42.24	32.21	39.26	71.20	59.70
AraPro	44.18	46.73	50.81	39.87	48.39	73.53	61.82
AraTruthfulQA	17.92	17.60	21.67	15.46	20.01	53.56	49.54
AraMath	5.63	5.19	7.26	5.61	6.41	26.35	38.68
AraEval	<b>72.02</b>	<b>75.37</b>	<b>77.67</b>	68.26	75.96	<b>82.66</b>	66.38
OpenAI Arabic MMMLU	71.33	74.69	75.89	<b>73.08</b>	<b>79.54</b>	80.20	64.45
Arabic MMLU	61.60	63.02	68.17	57.04	65.60	79.95	<b>69.25</b>
Vocabulary Token Statistics							
Arabic tokens	29,552	36,028	37,195	43,857	32,046	3,990	3,769
Arabic and math tokens	29,643	36,065	37,236	44,947	32,137	4,311	4,995

\*Tokenizer identical to AceGPT-V2 8B/70B's.

\*\*Tokenizer identical to AceGPT-V2 32B's.

Table 16: Vocabulary coverage across Arabic benchmarks and model tokenizers.